# Development of an IP Core Network Model for performance analysis of Third Generation Radio Access Networks
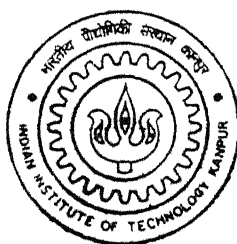
by

Mirza Shahrukh Baig

DEPARTMENT OF ELECTRICAL ENGINEERING

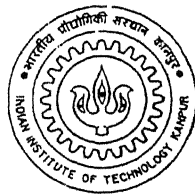INDIAN INSTITUTE OF TECHNOLOGY, KANPUR

March, 2001

# Development of an IP Core Network Model for performance analysis of Third Generation Radio Access Networks

*A Thesis Submitted*
*in Partial Fulfillment of the Requirements*
*for the Degree of*

*Master of Technology*

*by*

## Mirza Shahrukh Baig

*to the*

## Department of Electrical Engineering
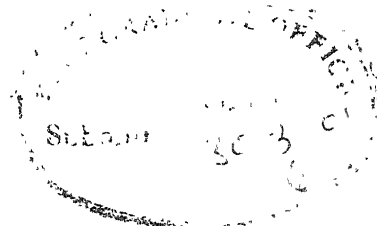## Indian Institute of Technology, Kanpur

## March, 2001

# Certificate

This is to certify that the work contained in the thesis entitled "*Development of an IP Core Network Model for performance analysis of Third Generation Radio Access Networks*", by *Mirza Shahrukh Baig*, has been carried out under my supervision along with the joint supervision of Prof. dr. Bernhard H. Walke of the University of Technology, Aachen, Germany and that this work has not been submitted elsewhere for a degree.

March, 2001

(Vishwanath Sinha)
Professor,
Department of Electrical Engineering,
Indian Institute of Technology,
Kanpur.

# Acknowledgement

I would like to express deep sense of gratitute to my thesis supervisor *Prof. Vishwanath Sinha* for his guidance and encouragement throughout the period of this thesis work. I hold him in reverential awe.

My sincere thanks are to *Prof. Bernhard Walke* in the chair of communication networks, RWTH-Aachen, Germany for providing me an opportunity to work at his institute. This was possible because of the Sandwich Exchange Programme between Technical Universities of Germany and IIT's, organized by the German Academic Exchange Service (DAAD). I would also like to thank *Dipl. Ing. Peter Stuckman* for his guidance during this thesis.

Mirza Shahrukh Baig

Dedicated to my parents

# ABSTRACT

With the fast evolution of *Global System for Mobile Communication* (GSM) towards third-generation (3G) mobile communication systems like *General Packet Radio Service* (GPRS)/*Universal Mobile Telecommunication Standard* (UMTS) new standards have to be integrated into the existing mobile radio networks. The driving force for this development is the predicted user demand for mobile data services that offer mobile Multimedia access and mobile Internet access. Since radio resources are scarce, *Quality of Service* (QoS) issues become very important for scalable use of these resources. Moreover as we move from circuit-switched to packet-switched services many changes have to be brought about in the existing networks to fulfill ever growing user demands like *voice, video, time-sensitive financial transactions, still images, large data files* and so on.

For performance analysis of 3G mobile data services both radio network and the core network have to be modeled to regard end-to-end *Quality of Service* behavior. Thus in this thesis a router model is developed which emulates the scheduling functions of an IP router as part of 3G Core Network based on *Third Generation Partnership Project* (3GPP) standards. Further an IP core network model (based on DiffServ) is set up which is then integrated into the ComNets GPRS simulator GPRSim. On the basis of simulations, this integration is evaluated for differentiated service performance. Finally DiffServ capable IP is compared with the present IP which is best effort.

# CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| | | | |
|---|---|---|---|
| AF | Assured Forwarding | CS | Coding Scheme |
| AMPS | Analog Mobile Phone Systems | DiffServ | Differentiated Services |
| ATM | Asynchronous Transfer Mode | DS | DiffServ |
| BA | Behavior Aggregate | DSCP | DiffServ Code Point |
| BCCH | Broadcast Control Channel | EDGE | Enhanced Data Rates for GSM Evolution |
| BLEP | Block Error Probability | EF | Expedited Forwarding |
| BS | Base Station | EGPRS | Enhanced General Packet Radio Service |
| BSC | Base Station Controller | | |
| BSS | Base Station Subsystem | ETSI | European Telecommunications Standards Institute |
| BSSGP | Base Station Subsystem GPRS Protocol | | |
| BTS | Base Transceiver Station | FDMA | Frequency Division Multiple Access |
| CAC | Connection Admission Control | | |
| CCCH | Common Control Channel | FIFO | First In First Out |
| CCH | Control Channel | FTP | File Transfer Protocol |
| CDMA | Code Division Multiple Access | GERAN | GSM/EDGE Radio Access Networks |
| C/I | Carrier to Interference Ratio | GGSN | Gateway GPRS Support Node |
| CN | Core Network | | |
| CNCL | Communication Networks Class Library | GMSC | Gateway Mobile Services Switching Centre |
| ComNets | Institute of Communication Netetworks | GPRS | General Packet Radio Service |
| COPS | Common Open Policy System | GPRSim | GPRS Simulator |
| | | GR | GPRS Register |

| | | | |
|---|---|---|---|
| GSM | Global System for Mobile Communication | PDCH | Packet Data Channel |
| | | PDN | Public Data Network |
| GSN | GPRS Support Node | PDP | Packet Data Protocol |
| GTP | GPRS Tunnelling Protocol | PDU | Protocol Data Unit |
| HLR | Home Location Register | PHP | Per Hop Behavior |
| HSCSD | High Speed Circuit Switched Data | PLMN | Public Land Mobile Network |
| HTML | Hypertext Markup Language | PQ | Priority Queuing |
| | | PSTN | Public Switched Telephone Network |
| HTTP | Hypertext Transfer Protocol | | |
| IETF | Internet Engineering Task Force | PTM | Point To Multipoint |
| | | PTP | Point To Point |
| IMSI | International Mobile Subscriber Identity | PDU | Protocol Data Unit |
| | | QoS | Quality of Service |
| IntServ | Integrated Services | RED | Random Early Detection |
| IP | Internet Protocol | RFC | Request For Comments |
| IPv4 | Internet Protocol Version 4 | RLC | Radio Link Control |
| IPv6 | Internet Protocol Version 6 | RLP | Radio Link Protocol |
| ISP | Internet Service Provider | RSVP | Resource ReSerVation Protocol |
| IS-95 | Interim Standard 95 | | |
| ITU | International Telecommunication Union | SDL | Specification and Description Language |
| LDP | Label Distribution Protocol | SDT | SDL Design Tool |
| LLC | Logical Link Control | SDU | Service Data Unit |
| MAC | Medium Access Control | SGSN | Serving GPRS Support Node |
| MF | Multi Field | | |
| MPLS | Multi Protocol Label Switching | SLA | Service Level Agreement |
| | | SMTP | Simple Mail Transfer Protocol |
| MS | Mobile Station | | |
| MSC | Mobile Services Switching Centre | SNMP | Simple Network Management Protocol |
| MSP | Multiple Subscriber Profile | SNDCP | Subnetwork Dependent Convergence Protocol |
| MT | Mobile Terminal | | |

| | | | |
|---|---|---|---|
| TACS | Total Access Communication Systems | UTRAN | UMTS Terrestrial Radio Access Networks |
| TCA | Traffic Conditioning Agreement | VLR | Visitor Location Register |
| TCH | Traffic Channel | VoIP | Voice over IP |
| TCP | Transmission Control Protocol | WAN | Wide Area Network |
| TDMA | Time Division Multiple Access | WAP | Wireless Application Protocol |
| TE | Terminal Equipment | WCDMA | Wideband Code Division Multiple Access |
| TFT | Traffic Flow Template | WFQ | Weighted Fair Queuing |
| TLLI | Temporary Logical Link Identity | WRC | World Radio Communication conference |
| TRX | Transceiver | WRED | Weighted Random Early Detection |
| UDP | User Datagram Protocol | WRR | Weighted Round Robin |
| UMTS | Universal Mobile Telecommunication System | WWW | World Wide Web |
| | | 3GPP | 3rd Generation Partnership Project |

# CHAPTER 1

# Introduction

Wireless data services are expected to see the same explosive growth in demand that wireline Internet services and wireless voice services have seen in recent years. Third Generation (3G) mobile devices and services will transform wireless communication into on-line, real time connectivity. 3G wireless technology will allow an individual to have immediate excess to location specific services that offer information on demand. The first generation of mobile phones consisted of analog models that emerged in early 1980s. The second generation of digital mobile phones appeared about ten years later along with first digital mobile networks. During the second generation, the mobile telecommunication industry experienced exponential growth both in terms of subscribers as well as new types of value added services. Mobile phones are rapidly becoming the preferred means of personal communication, creating the world's largest consumer electronics industry. The rapid and efficient deployment of new wireless data and Internet services has emerged as a critical priority for communications equipment manufacturers. Network components that enable wireless data services are fundamental to the next generation network infrastructure.

In the framework of evolution of the *Global System for Mobile communication* (GSM) towards third-generation (3G) mobile communication systems known as the *Universal Mobile Telecommunication System* (UMTS) or the *3rd Generation Partnership Project* (3GPP) new standards are presently integrated into existing mobile radio networks. Emerging Technologies like *High Speed Circuit Switched Data* (HSCSD) is already introduced in some countries. *The General Packet Radio Service* (GPRS) will be available in the year 2001 in Europe and many countries world wide. The next in line are EDGE (*Enhanced Data rates for GSM Evolution*) enhancements of GPRS and the entering of new air interfaces like *Wide-band Code*

*Division Multiple Access* (WCDMA) realizing *UMTS Terrestrial Radio Access Networks* (UTRANs).

For the interconnection of these new technologies with the information infrastructure, e.g. IP server farms or the public Internet all these radio access technologies will be based on the same core network architecture. Core networks standardized by *European Telecommunications Standards Institute* (ETSI)/*Third Generation Partnership Project* (3GPP) are composed of IP routers, which realize transmission of user data to gateway IP routers and the inter-working functions with subnetworks like external networks or other Public Land Mobile Networks (PLMNs). For performance analysis of 3G mobile data services both the radio network and the core network have to be emulated or modeled to regard end-to-end QoS behavior. The current Internet Protocol (IP) is not-sufficient for time sensitive traffic with specific latency, variance, packet loss, and bandwidth requirements. As a result, there is significant amount of ongoing research to improve the best effort services of IP. In the area of cellular communications, current second generation wireless systems cannot efficiently support data traffic with bursty, high bandwidth requirements. The development of 3G wireless networks is likewise producing a large volume of research. The combined advances in both these fields appear to promise QoS enhancements for future wireless and wired environments.

In the framework of this thesis an IP Core network model is developed which is then integrated into GPRS simulator for carrying out performance analysis in presence of both the radio access network and the core network. This IP Core network model is based on Differentiated Services architecture. This model consists of edge routers and a cascade of interior routers. All complex classification and conditioning is performed by edge routers with simple scheduling and queuing given to interior routers. In both edge routers and interior routers the scheduling algorithm that is used is weighted round robin (WRR) scheduling. The reason for choosing it was the differentiated services standard which says that the interior routers should be free from complex methods and work performed by interior routers should be a minimum. The results that came out clearly show the advantage of using differentiated services for mobile core networks as compared to the best effort IP.

In this thesis after the introduction, Chapter 2 provides an overview of Quality of Service aspects and QoS enhanced IP technologies like *Integrated Services* (IntServ), *Differentiated Services* (DiffServ) and the *Multi Protocol Label Switching* (MPLS).

Chapter 3 describes the differentiated services architecture and services in detail. This technology is seen at present to be very promising when it comes to scalability of network management resource. It is also the basis of the implementation of an IP core network model — the framework of this thesis.

Chapter 4 gives an overview of present and evolving mobile technologies. This is followed by Quality of Service issues in Third Generation Radio Access Networks.

Chapter 5 introduces the ComNets GPRS simulator GPRSim and Chapter 6 the IP Core Network simulator, and their integration.

Performance evaluation, discussion and interpretation of simulation results are given in Chapter 7.

Finally, we have conclusions in Chapter 8 followed by an outlook for future research in this fast growing field.

In this thesis after the introduction, Chapter 2 provides an overview of Quality of Service aspects and QoS enhanced IP technologies like *Integrated Services* (IntServ), *Differentiated Services* (DiffServ) and the *Multi Protocol Label Switching* (MPLS).

Chapter 3 describes the differentiated services architecture and services in detail. This technology is seen at present to be very promising when it comes to scalability of network management resource. It is also the basis of the implementation of an IP core network model – the framework of this thesis.

Chapter 4 gives an overview of present and evolving mobile technologies. This is followed by Quality of Service issues in Third Generation Radio Access Networks.

Chapter 5 introduces the ComNets GPRS simulator GPRSim and Chapter 6 the IP Core Network simulator, and their integration.

Performance evaluation, discussion and interpretation of simulation results are given in Chapter 7.

Finally, we have conclusions in Chapter 8 followed by an outlook for future research in this fast growing field.

# Quality of Service in IP Networks

## 2.1 QoS Definitions

In the simplest sense. Quality of Service (QoS) means providing a consistent. predictable data delivery service, in other words satisfying customer application requirements.

QoS refers to the ability of a network element (e.g. an application. host or router) to have some level of assurance that its traffic and service requirements can be satisfied. Naturally. QoS would require the cooperation of all network layers from top-to-bottom. as well as every network element from end-to-end. Any QoS assurances are only as good as the weakest link service assurance in the chain between a sender and a receiver [Sta99b]. QoS does not create bandwidth. It isn't possible for the network to give what it doesn't have, so bandwidth availability is a starting point. QoS only manages bandwidth according to application demands and network management settings, and in that regard it cannot provide certainty if it involves sharing. Hence, QoS with a guaranteed service level requires resource allocation to individual data streams. A priority for QoS designers has been to ensure that best-effort traffic is not starved after reservations are made. At the same time, QoS-enabled high-priority applications must not disable the low-priority Internet applications.

In broad terms the Quality of Service of a network is a measure of how well it does its job i.e. how quickly and reliably it transfers various kinds of data, including digitalized voice and video traffic from source to destination.

### 2.1.1 Need of Quality of service

Back when networks dealt pretty much exclusively with voice, the QoS hardly came up. The circuit switched telephone system was designed specifically to satisfy the

human ear. It did and it does. Today, with the advent of packet switching and entering of many kinds of communications traffic *(time-sensitive financial transactions. voice. video . still images, large data files. and so on)* there are more than one set of criteria to satisfy. The data rate needed for satisfactory voice communication may take an intolerable time to transfer high resolution images. Conversely. the degree of network latency acceptable in transferring some files may not be adequate for real-time voice. So QoS has become a hot topic as different applications require different QoS.

The Internet Protocol (IP), and the architecture of the Internet itself, is based on the simple concept that datagrams with source and destination addresses can traverse a network of (IP) routers independently, without the help of their sender or receiver. The Internet was historically built on the concept of a dumb network. with smarts at either end (at the sender and receiver).

There is a price to pay for this simplicity, however. The reason IP is simple is because it doesn't provide many services. IP provides addressing. and that enables the independence of each datagram. IP can fragment datagrams (in routers) and reassemble them (at the receiver). allowing traversal of different network media. But IP does not provide reliable data delivery. Routers are allowed to discard IP datagrams en route, without notice to sender or receiver. IP relies on upper-level transports (e.g. TCP) to keep track of datagrams, and retransmit as necessary. And these reliability mechanisms can only assure data delivery; neither IP nor its high-level protocols can ensure timely delivery or provide any guarantees about data throughput. IP provides what is called a best effort service. It can make no guarantees about when data will arrive, or how much it can deliver.

This limitation has not been a problem for traditional Internet applications like web, email, file transfer, and the like. But the new breed of applications, including audio and video streaming, demand high data throughput capacity (bandwidth) and have low-latency requirements when used in two-way communications (i.e. conferencing and telephony). Public and private IP networks are also being used increasingly for delivery of mission-critical information that cannot tolerate unpredictable losses.

Unlike pure virtual circuit technologies like ATM and Frame Relay, IP does not

available bandwidth, and it is also more flexible. Typical network traffic is bursty rather than continuous. IP is datagram-based so it uses the available bandwidth most efficiently, by sharing what is available as needed. This also allows IP to adapt more flexibly to applications with varying needs. However, it also leads to some unpredictability in service. The capacity to tolerate this unpredictability relates to the level of guarantee they require [Sta99b].

## 2.1.2 Quality of Service parameters

Technically, QoS refers to an aggregation of system performance metrics. The five most important of these are [DR00, Tri98]:

**Availability:** The availability of a network, its components, or even a service, should ideally approximate 100%. Even a high-sounding figure of as much as 99.8% translates into about an hour and a half of down time per month. This might not be acceptable to a large enterprise. Thus, consumer satisfaction largely depends on this parameter reaching as high as possible.

**Throughput:** This is the effective data transfer rate measured in bits per second. It depends on – but is explicitly not the same as – the maximum capacity or bandwidth of the network. Sharing a network lowers the throughput realizable by any user, as does the overhead imposed by the extra bits included in every packet for identification and other purposes. A minimum rate of throughput is usually guaranteed by a service provider.

**Packet loss:** Network elements, like switches and routers, are equipped with buffered queues to adopt to link congestion to some extent. However, if a link remains congested for too long, it may result in a buffer overflow and thus a loss of data. In a mobile radio network packets may additionally get lost owing to the special conditions on the radio interface. Both cases usually necessitates in a retransmission of the packet becoming necessary, increasing the total transmission time. In a well managed network, packet loss will typically be less than 1 percent averaged over, say, a month.

**Latency:** Latency is understood as the time taken by data to travel from its source to its destination. Thus, it may also be referred to as end-to-end delay. Unless satellites are involved, the latency of a 5000 km voice call carried by a circuit

switched telephone network is about 25 ms. For the present public Internet, a voice call may easily exceed 150 ms of latency because of delays, such as those caused by signal processing and congestion.

Various components add up to the end-to-end delay experienced by a packet on a transmission path:

- *Transmission delay:* the time it takes to put all bits of a packet onto the link

- *Propagation delay:* the time it takes for a bit to traverse a link (usually, at the speed of light)

- *Processing delay:* the time it takes to process a packet in a network element (e.g., routing it to the output port)

- *Queuing delay:* the time a packet must wait in a queue before it is scheduled for transmission

**Jitter:** Jitter, or latency variation, may be induced by various causes, e.g., variations in queue length, variations in the processing time needed to reorder packets that arrived out of order because they traveled over different paths, and variations in the processing time needed for re-assembly of packets segmented before being transmitted. This distortion is particularly damaging to multimedia traffic. For example, the playback of audio or video data may have a jittery or shaky quality.

Different applications, naturally, require different QoS; meaning that one application can be sensitive to one of the parameters while it might not be sensitive to another (see Figure 2.1). Voice over IP is sensitive to latency and if the level of latency reaches a certain level then the application becomes useless, on the other hand it is not very concerned with reliability (relatively). If a few bits are lost over a voice transmission it is not critical. But for a transmission of an executable program just a few bit errors might be devastating. It is these differences that must be considered in writing the Service Level Agreements (SLAs) between the service providers and their clients. The usual agreement specifies the end-to-end performance to which the client is entitled over a specified time interval.

Table 2.1: Varied sensitivities of network data types [DR00]

| Traffic type | Sensitivities | | | |
| --- | --- | --- | --- | --- |
| | Bandwidth | Loss | Delay | Jitter |
| Voice | Very low | Medium | High | High |
| E-commerce | Low | High | High | Low |
| Transactions | Low | High | High | Low |
| E-mail | Low | High | Low | Low |
| Telnet | Low | High | Medium | Low |
| Casual browsing | Low | Medium | Medium | Low |
| Serious browsing | Medium | High | High | Low |
| File transfers | High | Medium | Low | Low |
| Video conferencing | High | Medium | High | High |
| Multicasting | High | High | High | High |

Network applications can be characterized in terms of how predictable the data rate is (see Table 2.2), and how tolerant of delay delivery is (see Table 2.3). Generally, two way applications are more sensitive to delay than one-way applications.

## 2.2 Key QoS mechanisms

**Admission control:** Admission Control determines whether a requested connection is allowed to be carried by the network. The main considerations behind this decision are current traffic load, current QoS, requested traffic profile, requested QoS, pricing and other policy considerations.

**Traffic shaping/conditioning:** In QoS enabled IP networks, it's necessary to specify the traffic profile for a connection to decide how to allocate various network resources. Traffic shaping/conditioning ensures that traffic entering at an edge or a core node adheres to the profile specified.

**Packet classification and marking:** In order to provide the requested QoS, it's critical to classify packets to enable different QoS treatment. This can be done based on various fields in IP headers (e.g., source/destination addresses and protocol type) and higher layer protocol headers (e.g., source/destination port

Table 2.2: Terms to characterize application data rates in terms of relative predictability [Sta99b]

| Rate Type | Description |
| --- | --- |
| Stream | Predictable delivery at a relatively constant bit rate(CBR). For example, although their rates often fluctuate, audio and video data streams are considered CBR because they have a quantifiable upper bound. |
| Burst | Unpredictable delivery of blocks of data at a variable bit rate(VBR). Applications like file transfer move data in bulk that can increase data rate to use all available bandwidth and has no upper bound. |

Table 2.3: Terms to characterize application sensitivity to data delivery delays [Sta99b]

| Delivery Type | Description |
| --- | --- |
| Asynchronous | No constraints on delivery time |
| Synchronous | Data is time-sensitive, but flexible. |
| Interactive | Delays may be noticeable to users/applications, but do not adversely affects usability. |
| Isochronous | Time-sensitive to an extent that adversely affects usability. |
| Mission-Critical | Data delivery delays disable functionality. |

numbers for Transmission Control Protocol or User Datagram Protocol). Efficient and consistent packet classification is a key problem under active research. In Differentiated services packets are classified based on IP header's Type of Service (TOS) byte for IPv4 (Internet Protocol version 4) and Traffic Class byte for IPv6 (Internet Protocol version 6).

**Priority and scheduling mechanisms:**  To satisfy the QoS needs of different connections, nodes need to have priority and scheduling mechanisms. The priority feature typically refers to the capability of providing different delay treatment. e.g., higher priority packets are always served before the lower priority ones, both in the context of packet processing and transmission on outbound links.

Nodes also implement different loss priority treatment. i.e., higher loss priority packets are lost less often than the lower loss priority ones.

**Signaling protocols:** To obtain the required QoS from a network, end-systems need to signal the network the desired QoS as well as the anticipated offered traffic profile. This has been a fundamental part of various connection-oriented networks (e.g., ATM). However, for connectionless networks (e.g., IP), this is relatively new. Corresponding examples are the signaling associated with Resource ReSerVation Protocol (RSVP) and Label Distribution Protocol (LDP). Implementation scalability and the corresponding capabilities to signal different QoS needs are issues under current examination.

**Congestion Control:** For QoS IP networks to operate in a stable and efficient fashion, it's essential that they have viable and robust Congestion Control capabilities. These capabilities refer to the ability to flow control and shed excessive traffic during the periods of congestion e.g. Random Early Detection (RED). RED prescribes discard probability to drop packets in a fair and robust way (i.e., consistent with behavior of higher layer protocols like TCP) based on the measured average queue length. RED (Random Early Detection) attempts to avoid congestion rather than reacting to it (and thereby avoid TCP synchronization problems that can result when hosts decrease or increase TCP traffic simultaneously after congestion occurs). It randomly drops packets before queues fill, to keep them from overflowing.

## 2.3  Service Level Agreement(SLA)

**SLA Parameters:** A Service Level Agreement (SLA) is a service contract between Service provider and their customer that defines provider responsibilities in terms of network levels (**throughput, loss rate, delays and jitter**) and times of availability, method of measurement, consequences if service levels aren't met or the defined traffic levels are exceeded by the customer, and all costs involved. It specifies the forwarding service a customer should receive.

## 2.4  Policy

Quality of Service protocols provide the mechanics to differentiate traffic, and policy defines how they are used. In the simplest sense, policy is one or more rules that describe the action(s) to occur when specific condition(s) exist. A service is the expression of a relationship between a set of objects whereas policy is a statement about a set of relationships between objects that provide a particular service. A policy can be used to configure and control a service. In an open and public Internet (as well as large intranets), the acceptance of QoS requests results in better network service to some flows, possibly at the expense of service to traditional best-effort flows. The purpose of such classification may include preferential queuing or dropping, admitting or denying access, or encrypting the packet's payload, to cite just a few examples. Protocols that explicitly support some or all of these functions include COPS (common open policy system), RSVP, IntServ, and DiffServ etc. The multiple types of devices that must work in concert across even a single domain to achieve the desired policy can include hosts (clients and servers), routers, switches, firewalls, bandwidth brokers, subnet bandwidth managers, network access servers, and policy servers [Sta99a]. The *Internet Engineering Task Force* (IETF) Policy working group is working on how to represent, manage, and share policies and policy information in a vendor-independent, inter-operable, scalable manner for QoS traffic management.

## 2.5  Importance of Priority

QoS is largely about priorities. At network aggregation points, like routers, multiplexers, and switches, data streams with different QoS needs are combined for transport over a common infrastructure. Satisfactory QoS has two main requirements: a means for labeling flows with respect to their priorities, and network mechanisms for recognizing the labels and acting on them. Some networks, notably those that use the asynchronous transfer mode (ATM) have extensive provisions of this kind. Unfortunately, the Internet does not have the same, and neither do the similar IP networks based on the transmission control protocol/ Internet protocol (TCP/IP)

suite. So ensuring adequate QoS comes down to devising a means for labeling data flows and recognizing and acting on those labels.

IP is a best-effort protocol in that it does not guarantee delivery of data packets. Confirmation of the arrival of data packets at the destination is the responsibility of the TCP. If any packet is not delivered (as determined by checking the sequence numbers of packets at the destination), TCP requests a retransmission of the missing packet. thereby ensuring that all packets eventually get to the destination. This is effective, but slow. Therefore, TCP is generally used by applications that are not time-sensitive.

Real-time applications cannot take advantage of TCP. Obviously. the time needed for keeping track of missing packets and retransmitting them is not acceptable in such cases. So these applications rely on what is essentially a stripped-down version of TCP, known as the user datagram protocol (UDP), which runs faster than TCP by omitting some of its functionality. Applications that run over UDP must either have those missing capabilities built into them or else do without.

In the case of voice communications, where retransmitting packets takes too long to be of any value anyway, missing packets are simply lost. Internet telephony. therefore. will work only over networks that are quite reliable to begin with, like fiber-based nets with relevent switches and routers.

## 2.6   QoS Architecture Models for Traffic Engineering

### 2.6.1   Integrated Services (IntServ)

The IntServ model modifies the basic IP service model to provide QoS based upon end-to-end resource reservation. IntServ provides QoS by allocating network resources to individual packet flows. For this purpose a flow is defined as a stream of IP packets between applications on end hosts, which have the same source and destination addresses, (TCP/UDP) port numbers, and protocol field. The performance requirements of a specific flow are known as the flowspec, and can include requirements such as bandwidth and delay. Unlike a *circuit switched* voice network where each call has the same bandwidth and delay requirements, each flow in the IntServ model can have a unique flowspec. The IntServ model is intended to allow

both time-sensitive and non-time-sensitive traffic flows to be serviced by the same IP layer.

Three main components of the IntServ model are a signaling protocol. an admission control capability, and a packet forwarding mechanism. IntServ uses the Resource Reservation Protocol (RSVP) to allocate network resources to each flow. If RSVP can not reserve sufficient network resources to meet the flowspec. then admission control is used to prevent the flow from entering the network. In addition IntServ requires a packet forwarding mechanism that can classify packets. manage buffers. and schedule service in the network routers. Thus in order to implement IntServ. a router must be able to maintain flow state information for each traffic flow that it services. The IntServ model suffers from scalability issues based upon providing QoS on a per flow basis. IntServ requires a router to maintain state information for all flows that pass through it. This becomes a significant problem as link speeds increase to gigabit-per-second or tetrabit-per-second rates with a corresponding increase in the number of simultaneous flows. In addition. every router along a path must use the IntServ model of traffic service in order to provide end-to-end resource reservation and QoS packet forwarding. While it is possible to upgrade all routers within a small corporate network to provide IntServ based QoS, this can not be easily accomplished in the Internet. DiffServ is a newer approach to providing QoS enhancements to IP that addresses the scalability problem.

## 2.6.2  Differentiated services (DiffServ)

At present DiffServ is seen as very promising as far as scalability is concerned. In DiffServ, 8 bit field present in the IP header called the DS (DiffServ) field is used to classify packets. Data flows having the same resource requirements may then be aggregated on the basis of their DS field when they arrive at the edge routers. The routers at the core can then forward the data flows toward their destinations depending upon the DiffServ field. Since most of the decision-making is in this way transferred from the core routers to the edge routers, the core network runs much faster.

In the past. QoS planners supported both IntServ and DiffServ. At present. however. the trend is to use DiffServ supplemented by some of the resource reservation capabilities of RSVP.

As this is the basis of the IP core network model in the framework of this thesis is explained in detail in the next chapter.

## 2.6.3  Multi Protocol Label Switching (MPLS)

It is another approach to speeding the transit of data through a network. Like IntServ and DiffServ it is also promulgated by the IETF.

Normally, under IP, packet headers are examined at every transit point (multiplexer. router, or switch) in a network, which takes time and contributes to the overall data delay. A more efficient approach would be to label the packets in such a way as to make it unnecessary for each IP packet header to be analyzed at points intermediate between the source and destination. Multiprotocol label switching does this by appropriately labeling IP packets at the input of label edge routers located at the entry points of an MPLS-enabled network.

The procedure works like this: the label edge router examines the incoming packets and decides based on the packet's source address, destination address, and priority level where to send it for its next hop through the network. It also attaches a 32-bit tag, known as an MPLS label, to the packet. The MPLS label contains such information as whether the packet should be treated as MPLS traffic or routed as an ordinary IP packet; whether it conforms to IPv4 or IPv6; the packet's time to live; and, of course, what its next hop should be. The edge router then forwards the packet to the router at the end of the next hop.

That router, in turn, examines the MPLS label and decides on the next hop for the packet. That second router then creates a second MPLS label. The two labels are swapped before the packet is forwarded to the second hop. The process is repeated until the packet reaches its destination.

This procedure has two advantages over normal IP routing. First of all, the routers along the path need not read and analyze a packet's complete header information, just the shorter MPLS label. This alone saves some time. Secondly, the swapping of labels leaves a trail in the registry of the routers that other packets in the same

session can follow. Once the first packet establishes a path, decision-making at intermediate points is eliminated to a great extent. This markedly speeds up the transfer of data [DR00].

### 2.6.4   MPLS with DiffServ

A future approach may be to use both MPLS and DiffServ to improve performance. Even at present some research is going on in this field. MPLS specifies ways to map Layer 3 traffic to a connection-oriented Layer 2 transports like ATM and frame relay. It adds a label containing specific routing information and allows routers to assign an explicit path to various classes of traffic. The MPLS network provides a transit service to DiffServ traffic by mapping DS-field to an explicit path through the MPLS network. For each Service Class, MPLS network provides an independent data flow path between the edge routers. The edge router obtains QoS information from the DS-field contained in the packet arriving from a DiffServ domain, and this information is used to map to an explicit path through the MPLS network.

The similarities and differences among the various procedures for controlling quality of service (QoS) are best understood by viewing them in the same general framework as shown in Figure 2.1. Data from one or more applications pass down through QoS enablers, which prioritize the data flows and indicate the resources each requires. The data then continues through various levels of software and hardware that control packet discard mechanisms when buffered queues becomes too long. Finally it reaches the basic transport mechanisms and their hardware platforms that carry packet to the next node.

## 2.7   Scope of QoS service and QoS Domains

If QoS is to be provided end-to-end between two end users, the scope of QoS becomes quite important. The term scope refers to the topological extent over which the customer is given the QoS [Bla00], (see Figure 2.2). The scope may be restricted to one provider, or many. But for meaningful QoS to be obtained, the scope of service should include all providers that are involved in handling the user's traffic. The concept has become more visible in the past few years as customers have had their traffic transported through one or more than one service provider.

| Applications | | | | | |
|---|---|---|---|---|---|
| Voice over IP (VoIP) | File transfer protocol (FTP) | E-mail | Mission-critical data | Audio, video, and teleconferencing | Virtual private network (VPN) |

| QoS enablers (classification of packets, prioritization, and indication of resources solicited) | | | |
|---|---|---|---|
| INTSERV | DIFFSERV | MPLS | Other present or future enablers |
| Resource reservation protocol (RSVP) | | DiffServ code point (DSCP) | Other present or future resource reservation or prioritizing schemes |

| Congestion avoidance and packet discard |
|---|
| Congestion management (queuing) |
| Link efficiency management (compression and shaping of data for a smoother flow) |

| Basic transport mechanisms: ATM, frame relay, cable modem, digital subscriber line, and others |
|---|
| Hardware platforms: routers and switches |

Figure 2.1: Similarities and differences of various QoS procedures

As long as the end-to-end service was best effort (no attempt to provide QoS to the customer), the scope of service was not very important. Certainly, the customer needed some level of service end-to-end, but services such as delay, jitter, and throughput were moot points, since they were simply not available. Today, with the increased importance of QoS, the scope of service becomes a key part of a service level agreement (SLA).

The scope is defined between an ingress point (where user traffic enters a QoS node or network) and an egress point (where user traffic leaves the QoS node or network). If user traffic spans multiple service providers, it is important that these providers have SLA agreement among themselves in relation to the customer's QoS needs. This idea of multi-QoS domain is called a QoS region. The agreements for multidomain service are best met by providing quantitative services.

Figure 2.2: Scope of service and QoS domains

# CHAPTER 3

# Differentiated Services (DiffServ)

## 3.1    Introduction

Today, the aim of networks is to assign each application exactly as much resources as it requires. no less but no more. As *DiffServ* promises this criteria better than *IntServ*, it is very popular at present in the IETF and among the vendors. Another reason of its popularity is its simplicity. There is strong belief that in future gigabit networks, the QoS model has to be simple. Differentiated services grew out of desire to find an approach that would be simple, scalable and relatively easy to deploy in a predominantly best effort Internet. In addition, within differentiated services there is a significant emphasis on allowing for meaningful end-to-end services to be provisioned across multiple, separately administered network clouds and on keeping the consequent business models as simple as possible.

*Differentiated services* (DiffServ) can be defined as a novel approach to providing *Quality of Service* (QoS) within the Internet in a scalable manner. It is a relatively simple and coarse method of providing differentiated classes of service for Internet traffic, to support various applications, and specific business requirements. The differentiated services approach to providing quality of service in networks employs a small, well-defined set of building blocks from which a variety of aggregate behaviors may be built. A small bit-pattern in each packet, in the IPv4 TOS octet or the IPv6 Traffic Class octet, is used to mark a packet to receive a particular forwarding treatment, or per-hop behavior, at each network node. A common understanding about the use and interpretation of this bit-pattern is required for inter-domain use, multi-vendor interoperability, and consistent reasoning about expected aggregate behaviors in a network. Thus, the DiffServ Working Group (IETF) has standardized a common layout for a six-bit field of both octets, called the 'DS field' [NBBB98].

## 3.2  Differentiated Services Architecture

Differentiated Services (DiffServ) is a set of technologies which allow network service providers to offer services with different kinds of network quality-of-service (QoS) objectives to different customers and their traffic streams.

The premise of DiffServ networks is that routers within the core of the network handle packets in different traffic streams by forwarding them using different per-hop behaviors (PHBs). The PHB to be applied is indicated by a DiffServ codepoint (DSCP) in the IP header of each packet. The DSCP markings are applied either by a trusted customer or by the edge routers on entry to the DiffServ network. The advantage of such a scheme is that many traffic streams can be aggregated to one of a small number of behaviour aggregates (BA) which are each forwarded using the same per-hop behavior at the router, thereby simplifying the processing and associated storage. In addition, there is no signaling, other than what is carried in the DSCP of each packet, and no other related processing that is required in the core of the DiffServ network since QoS is invoked on a packet-by-packet basis.

This architecture achieves scalability by implementing complex classification and conditioning functions only at network boundary nodes, and by applying per-hop behaviors to aggregates of traffic which have been appropriately marked using the DS field in IPv4 and IPv6 headers. Per-application flow or per-customer forwarding state need not be maintained within the core of the network. Network resources are allocated to traffic streams by service provisioning policies which govern how traffic is marked and conditioned upon entry to a differentiated services-capable network, and how that traffic is forwarded within that network. Main features of Differentiated Services can be seen from Table 3.1.

The differentiated services architecture is based on a simple model where traffic entering a network is classified and possibly conditioned at the boundaries of the network, and assigned to different behavior aggregates (see Figure 3.1). Each behavior aggregate is identified by a single DS codepoint. Within the core of the network, packets are forwarded according to the per-hop behavior associated with the DS codepoint [BBC+98]. The architecture consists of a set of necessary functional elements, which are as follows:

Table 3.1: Differentiated Services (DiffServ) features

---

* Services differentiated by performance (and may be price)

* Service on packet-by-packet basis

* Does not define a control plane (a control or signaling protocol)

* Attempts to force complexity out of network to edges

* Concerned with: Traffic classification and Traffic conditioning

* Relies on IP header to contain a label (a codepoint) to identify traffic type

* Requires rules for traffic conditioning for metering, marking, shaping, policing

* Rules are called Traffic Conditioning Agreement (TCA)

---

**Per-Hop Behaviors at Interior Routers:** An interior router is any router not at the boundary of a differentiated services network domain. Since interior routers make up the vast majority of routers through which most IP packets pass, the complexity of the functions performed by interior routers must remain low. The differentiated services architecture recognizes this fact and mandates that only simple per-hop behaviors(PHBs) be implemented at interior routers. A per-hop behavior is any forwarding behavior performed by the router, and it usually consists of a packet queuing and scheduling policy. Interior routers select a PHB for a packet by examining its Differentiated Services (DS) field. which is contained within the IP header (Type of Service(TOS) octet of IPv4. or Traffic Class octet of IPv6).

**Traffic Classification and Conditioning at Boundary Routers:** As its name implies, a boundary router exists at the edge of a differentiated services network domain. This boundary must perform sophisticated packet classification, metering, marking, policing, and shaping operations of packets arriving at it.

## 3.2.1 Differentiated Services Domain

A DS domain is a contiguous set of DS nodes which operate with a common service provisioning policy and set of PHB groups implemented on each node (see Figure 3.2). A DS domain has a well-defined boundary consisting of DS boundary nodes which classify and possibly condition ingress traffic to ensure that packets which transit the domain are appropriately marked to select a PHB from one of the

Figure 3.1: Differentiated Service Architecture



Figure 3.2: DiffServ Domain [Bla00]

PHB groups supported within the domain. Nodes within the DS domain select the forwarding behavior for packets based on their DS codepoint. Inclusion of non-DS-compliant nodes within a DS domain may result in unpredictable performance and may impede the ability to satisfy service level agreements. A DiffServ domain normally consists of one or more networks under the same administration; for example, an organization's intranet or an ISP. The administration of the domain is responsible for ensuring that adequate resources are provisioned or reserved to support the SLAs offered by the domain.

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
|   |   |   |   |   |   |   |   |

DiffServ Code Point (DSCP) ◄─────────────────► Currently Unused (CU) ◄──────►

Figure 3.3: IPv4 Type of Service(TOS) octet or IPv6 Traffic Class octet

### 3.2.2  Differentiated Services Region

A differentiated services region (DS Region) is a set of one or more contiguous DS domains. DS regions are capable of supporting differentiated services along paths which span the domains within the region. The DS domains in a DS region may support different PHB groups internally and different codepoint to PHB mappings. However, to permit services which span across the domains, the peering DS domains must each establish a peering SLA which defines (either explicitly or implicitly) a TCA which specifies how transit traffic from one DS domain to another is conditioned at the boundary between the two DS domains. It is possible that several DS domains within a DS region may adopt a common service provisioning policy and may support a common set of PHB groups and codepoint mappings, thus eliminating. the need for traffic conditioning between those DS domains.

## 3.3  DiffServ Field (DS Field)

The DiffServ architecture is based on the use of the DS Field which is placed in the IPv4 Type of Service octet or IPv6 Traffic Class octet [NBBB98]. First Six bits of the DS field are used as a codepoint (DSCP) to select the PHB a packet experiences at each node. A two-bit currently unused (CU) field is reserved for future use (see Figure 3.3). The value of the CU bits are ignored by differentiated services-compliant nodes when determining the per-hop behavior to apply to a received packet. In a DSCP value notation 'xxxxxx' (where 'x' may equal '0' or '1') the left-most bit signifies zeroeth bit of the DS field, and the right-most bit signifies fifth bit. DS-compliant nodes must select PHBs by matching against the entire 6-bit DSCP field,

Table 3.2: DiffServ(DS) Field features

---

* Bits 0.1.2 : used for priority setting and further clarifies their functions

* Bits 3.4.5 · used for finer granularity through packet drop probability

* Bits 6,7 : currently unused

* Backward compatible with precedence priority(RFC791) by reserving code space xxx 000 xx for it

* Code 000 000 xx is reserved for best effort traffic

* Code point space xxx xx0 xx (32 points) is allocated to standard use

* Code point space xxx x11 xx (16 points) for local use

* Code point space xxx x01 xx is reserved for future use

---

e.g., by treating the value of the field as a table index which is used to select a particular packet handling mechanism which has been implemented in that device.

DiffServ field features can be seen from Table 3.2. The proposed DiffServ standard is backward compatibility with RFC 791 implementations (i.e. precedence priority), but allows more efficient use of the 3rd, 4th and 5th bits. The DiffServ standard utilizes the most significant bits 0, 1, and 2 for priority setting, but further clarifies their functions, plus offers finer priority granularity through use of the next three bits in the DiffServ field. DiffServ names the priority levels (defined by the three most significant bits of the DS Field) into the categories as shown in Table 3.3. Bits 3rd and 4th of the TOS field (now called the DSCP in the DiffServ standard) allow further priority granularity through the specification of a packet drop probability for any of the defined class as shown in Table 3.4. Bits 0, 1, and 2 define the class; bits 3 and 4 specify the drop percentage; bit 5 is always 0. Using this system, a device would first prioritize traffic by class, then differentiate and prioritize same-class traffic by considering the drop percentage. It is important to note that this standard has not specified a precise definition of low, medium, and high drop percentages. Additionally, not all devices will recognize the DiffServ bit 3 and 4 settings. Also even when the settings are recognized, they do not necessarily trigger the same forwarding action to be taken by each type of device on the network—each device will implement its own response in relation to the packet priorities it detects. The DiffServ proposal

Table 3.3: Precedence Levels of DiffServ based on bits 0, 1, 2 of DSCP

| 0,1,2 bits of DSCP | Precedence Level | |
|---|---|---|
| 111 | Precedence 7 | used for link layer and routing protocols |
| 110 | Precedence 6 | used for IP routing protocols |
| 101 | Precedence 5 | Express Forwarding (Expedited forwarding) |
| 100 | Precedence 4 | Assured Forwarding class 4 |
| 011 | Precedence 3 | Assured Forwarding class 3 |
| 010 | Precedence 2 | Assured Forwarding class 2 |
| 001 | Precedence 1 | Assured Forwarding class 1 |
| 000 | Precedence 0 | Best Effort Class |

Table 3.4: Drop Precedence Classification based on bits 3, 4, 5 of DSCP

| | Low Drop precedence | Medium Drop precedence | High Drop precedence |
|---|---|---|---|
| Class 1 (first six bits of DSCP) | 001010 | 001100 | 001110 |
| Class 2 (first six bits of DSCP) | 010010 | 010100 | 010110 |
| Class 3 (first six bits of DSCP) | 011010 | 011100 | 011110 |
| Class 4 (first six bits of DSCP) | 100010 | 100100 | 100110 |

is meant to allow a finer granularity of priority setting for the applications and devices that can make use of it, but it does not specify interpretation.

## 3.4   Traffic Classification and conditioning

Differentiated services are extended across a DS domain boundary by establishing a Service Level agreement between a network and a DS domain. The SLA may specify packet classification and re-marking rules and may also specify traffic profiles and actions to traffic streams which are in-profile or out-of-profile. The TCA between the domains is derived from this SLA.

Traffic conditioning performs metering, shaping, policing or re- marking to ensure that the traffic entering the DS domain conforms to the rules specified in the TCA, in accordance with the domain's service provisioning policy as shown in Figure 3.4.

Traffic Conditioning Block(TCB)

Meter

Packet stream → Classifier → Marker → Shaper/ Dropper →

Figure 3.4: The DiffServ traffic classification and conditioning model

The extent of traffic conditioning required is dependent on the specifics of the service offering, and may range from simple codepoint re-marking to complex policing and shaping operations.

### 3.4.1 Classifiers

Packet classifiers select packets in a traffic stream based on the content of some portion of the packet header. In DiffServ two types of classifiers are defined. The BA (Behavior Aggregate) Classifier classifies packets based on the DS codepoint only. The MF (Multi-Field) classifier selects packets based on the value of a combination of one or more header fields, such as source address, destination address, DS field, protocol ID, source port and destination port numbers, and other information such as incoming interface.

Classifiers are used to steer packets matching some specified rule to an element of a traffic conditioner for further processing. Classifiers must be configured by some management procedure in accordance with the appropriate TCA.

This Classification is performed by a classifier element. Classifiers are 1:N (fan-out) devices, they take a single traffic stream as input and generate N logically separate traffic streams as output. Classifiers are parameterized by filters and output streams. Packets from the input stream are sorted into various output streams by

Classified
Traffic



Figure 3.5: Classifier example

filters which match the contents of the packet. The Classifier shown in Figure 3.5 separates traffic into one of three output steams.

## 3.4.2   Traffic Conditioners

A traffic conditioner contains : meter, marker, shaper, and dropper. A traffic stream is selected by a classifier, which steers the packets to a logical instance of a traffic conditioner. A meter is used (where appropriate) to measure the traffic stream against a traffic profile. The state of the meter with respect to a particular packet (e.g., whether it is in- or out- of-profile) may be used to affect a marking, dropping, or shaping action. A traffic conditioner may not necessarily contain all four elements. For example, in the case where no traffic profile is in effect, packets may only pass through a classifier and a marker.

**Meters:**   Traffic meters measure the temporal properties of the stream of packets selected by a classifier against a traffic profile specified in a TCA. A meter passes state information to other conditioning functions to trigger a particular action for each packet which is either in- or out-of-profile (to some extent).
A meter, measures the rate at which packets making up a stream of traffic pass it, compares the rate to some set of thresholds and produces some number of potential results: a given packet is said to be conformant to a level of the meter if, at the time that the packet is being examined, the stream appears to be within the rate limit for the profile associated with that level.

**Markers:**   Packet markers set the DS field of a packet to a particular codepoint, adding the marked packet to a particular DS behavior aggregate. The marker may be configured to mark all packets which are steered to it to a single code-

point, or may be configured to mark a packet to one of a set of codepoints used to select a PHB in a PHB group, according to the state of a meter. When the marker changes the codepoint in a packet it is said to have re-marked the packet.

**Shapers:** Shapers delay some or all of the packets in a traffic stream in order to bring the stream into compliance with a traffic profile. A shaper usually has a finite-size buffer, and packets may be discarded if there is not sufficient buffer space to hold the delayed packets.

**Droppers:** Droppers discard some or all of the packets in a traffic stream in order to bring the stream into compliance with a traffic profile. This process is known as policing the stream. Note that a dropper can be implemented as a special case of a shaper by setting the shaper buffer size to zero (or a few) packets.

## 3.5   Per Hop Behaviors (PHBs)

A per-hop behavior (PHB) is a description of the externally observable forwarding behavior of a DS node applied to a particular DS behavior aggregate. In the event that only one behavior aggregate occupies a link, the observable forwarding behavior (i.e., loss, delay, jitter) will often depend only on the relative loading of the link. Useful behavioral distinctions are mainly observed when multiple behavior aggregates compete for buffer and bandwidth resources on a node. The PHB is the means by which a node allocates resources to behavior aggregates, and it is on top of this basic hop-by-hop resource allocation mechanism that useful differentiated services may be constructed.

The most simple example of a PHB is one which guarantees a minimal bandwidth allocation of X percent of a link (over some reasonable time interval) to a behavior aggregate. This PHB can be fairly easily measured under a variety of competing traffic conditions. A slightly more complex PHB would guarantee a minimal bandwidth allocation of X percent of a link, with proportional fair sharing of any excess link capacity. In general, the observable behavior of a PHB may depend on certain constraints on the traffic characteristics of the associated behavior aggregate, or the characteristics of other behavior aggregates. PHBs may be specified in terms of their resource (e.g., buffer, bandwidth) priority relative to other PHBs, or in terms of their

relative observable traffic characteristics (e.g., delay, loss). PHB groups will usually share a common constraint applying to each PHB within the group, such as a packet scheduling or buffer management policy. The relationship between PHBs in a group may be in terms of absolute or relative priority. PHBs are implemented in nodes by means of some buffer management and packet scheduling mechanisms. PHBs are defined in terms of behavior characteristics relevant to service provisioning policies, and not in terms of particular implementation mechanisms. In general, a variety of implementation mechanisms may be suitable for implementing a particular PHB group. A PHB is selected at a node by a mapping of the DS codepoint in a received packet.

### 3.5.1  Expedited Forwarding (EF) PHB

The EF PHB can be used to build a low loss, low latency, low jitter, assured bandwidth, end-to-end service through DS domains. Such a service appears to the endpoints like a *point-to- point connection* or a *virtual leased line*. This service has also been described as premium service [JNP99]. Loss, latency and jitter are all due to the queues traffic experiences while transiting the network. Therefore providing low loss, latency and jitter for some traffic aggregate means ensuring that the aggregate sees no (or very small) queues. Queues arise when (short-term) traffic arrival rate exceeds departure rate at some node. Thus a service that ensures no queues for some aggregate is equivalent to bounding rates such that, at every transit node, the aggregate's maximum arrival rate is less than that aggregate's minimum departure rate (see Figure 3.6).

Thus the EF PHB is defined as a forwarding treatment for a particular DiffServ aggregate where the departure rate of the aggregate's packets from any DiffServ node must equal or exceed a configurable rate. Creating such a service has two parts:

1. Configuring nodes so that the aggregate has a well-defined minimum departure rate. (Well-defined means independent of the dynamic state of the node. In particular, independent of the intensity of other traffic at the node.)

2. Conditioning the aggregate (via policing and shaping) so that its arrival rate at any node is always less than that node's configured minimum departure rate.

Figure 3.6: Expedited Forwarding method

The EF PHB provides the first part of the service. The network boundary traffic conditioners provide the second part. Codepoint 101110 is recommended for the EF PHB.

## 3.5.2 Assured forwarding (AF) PHB

Assured Forwarding (AF) PHB group is a means for a provider DS domain to offer different levels of forwarding assurances for IP packets received from a customer DS domain. Four AF classes are defined, where each AF class is in each DS node allocated a certain amount of forwarding resources (buffer space and bandwidth) [HFB+99]. IP packets that wish to use the services provided by the AF PHB group are assigned by the customer or the provider DS domain into one or more of these AF classes according to the services that the customer has subscribed to.

Within each AF class IP packets are marked (again by the customer or the provider DS domain) with one of three possible drop precedence values. In case of congestion, the drop precedence of a packet determines the relative importance of the packet within the AF class. A congested DS node tries to protect packets with a lower drop precedence value from being lost by preferably discarding packets with a higher drop precedence value.

Table 3.5: Recommended values for AF codepoints

| AFij | Codepoint |
|------|-----------|
| AF11 | 001010 |
| AF12 | 001100 |
| AF13 | 001110 |
| AF21 | 010010 |
| AF22 | 010100 |
| AF23 | 010110 |
| AF31 | 011010 |
| AF32 | 011100 |
| AF33 | 011110 |
| AF41 | 100010 |
| AF42 | 100100 |
| AF43 | 100110 |

In a DS node, the level of forwarding assurance of an IP packet thus depends. on (1) how much forwarding resources has been allocated to the AF class that the packet belongs to, (2) what is the current load of the AF class, and, in case of congestion within the class, (3) what is the drop precedence of the packet.

Assured Forwarding (AF) PHB group provides forwarding of IP packets in N independent AF classes. Within each AF class, an IP packet is assigned one of M different levels of drop precedence. An IP packet that belongs to an AF class i and has drop precedence j is marked with the AF codepoint AFij, where $1 <= i <= N$ and $1 <= j <= M$. Currently, four classes (N=4) with three levels of drop precedence in each class (M=3) are defined for general use as shown in Table 3.5. More AF classes or levels of drop precedence may be defined for local use. Packets in one AF. class must be forwarded independently from packets in another AF class, i.e., a DS node must not aggregate two or more AF classes together.

## 3.6 Scheduling Mechanisms and Congestion Control

### 3.6.1 Scheduling Mechanisms

One way that network elements handle an overflow of arriving traffic is to use a queuing algorithm to sort the traffic, then determine some method of prioritizing it onto an output link. Some important scheduling mechanisms are [Cis00, Eri99]:

**First In, First Out (FIFO) Queuing:** In its simplest form, FIFO queuing involves storing packets when the network is congested and forwarding them in order of arrival when the network is no longer congested. FIFO is the default queuing algorithm in some instances, thus requiring no configuration. but it has several shortcomings. Most importantly, FIFO queuing makes no decision about packet priority: the order of arrival determines bandwidth. promptness. and buffer allocation. Nor does it provide protection against ill-behaved applications (sources). Bursty sources can cause high delays in delivering time-sensitive application traffic, and potentially to network control and signaling messages. FIFO queuing was a necessary first step in controlling network traffic, but today's intelligent networks need more sophisticated algorithms. As DiffServ is about different priority classes, FIFO cannot be used with Differentiated Services.

**Priority Queuing (PQ):** PQ ensures that important traffic gets the fastest handling at each point where it is used. It was designed to give strict priority to important traffic. Priority queuing can flexibly prioritize according to network protocol, incoming interface, packet size, source/destination address, and so on. In PQ, each packet is placed in one of N queues (say 4)—High, Medium. Normal, or Low—based on an assigned priority. During transmission, the algorithm gives higher-priority queues absolute preferential treatment over low-priority queues. This is a simple and intuitive approach but can cause queuing delays that the higher-priority traffic might have experienced to be randomly transferred to the lower-priority traffic, increasing jitter on the lower-priority traffic. PQ is useful for making sure that mission-critical traffic traversing various links gets priority treatment. A strict priority mechanism between two or more classes aims to provide the lowest possible delay for the highest priority

class. This mechanism sends the data from the highest priority class before sending data for the next class. This could lead to starvation of lower priority classes. Thus it is also not used with DiffServ.

**Weighted Round Robin (WRR) Queuing:** Weighted Round Robin (WRR) aims to give a weighted access to the available bandwidth to each class, ensuring a minimum allocation and distribution. The scheduling services each class in a round robin manner according to the weights. If one or more classes is not using its full allocation, then the unused capacity is distributed to the other classes according to their weighting. WRR ensures that queues do not starve for bandwidth and that the traffic gets predictable service. It serves each queue in a round robin manner, and for each turn, a number of bits corresponding to the queue's weight is pulled out from the queue. Weighted Round Robin is simple to implement. With the advent of technology, as we moving from links speed of Mbps to Gbps (giga bits per second)/Tbps(tetra bits per second) this may become a deciding factor, if the link speeds increase faster than the pure processing power.

**Weighted Fair Queuing (WFQ):** Weighted Fair Queuing (WFQ) like WRR aims to distribute available bandwidth over a number of weighted classes. The scheduling mechanism uses a combination of weighting and timing information to select which queue to service. The weighting effectively again controls the ration of bandwidth distribution between classes under congestion, and can also indirectly control delay for under utilized classes. WFQ could allow the unused capacity to be distributed differently from the minimum bandwidth weightings, such as a different configured weighting, or it could be dependent on traffic load in each class. Both WRR and WFQ can be used for scheduling mechanisms in DiffServ capable routers.

### 3.6.2 Congestion Control Tools

Rather simple congestion control mechanisms, typically based on discarding packets when buffers filled up to a threshold level, have been used in packet nodes. This would protect the routers, but gave rise to an unexpected phenomenon known as global synchronization. When TCP packets are discarded, the TCP scheduling

algorithm responds by lowering its transmission rate, then building it up again. When core routers overload, they drop packets from many hosts, leading to many TCP sessions backing off and ramping up their transmission rates synchronously again. This could lead to a saw tooth pattern of under-utilization and congestion [Eri99].

- Random early Detection(RED):

    To achieve higher average utilization, Random Early Detection (RED) was introduced to attempt to stop the synchronized effect. In addition to that the algorithm aims to provide fairness among the TCP flows competing for the resources. Instead of discarding all traffic when a threshold is reached. the RED algorithm start dropping incoming packets, selected by a random function. with increasing probability as buffer utilization increases towards the maximum. The aim is to make some proportion of TCP sessions back off before hitting congestion. As in DiffServ we have different classes and each of these class have three different drop precedence levels. Thus RED cannot be used effectively with Differentiated Services.

- Weighted Random Early Detection (WRED)

    WRED combines the capabilities of the RED algorithm with Differentiated Services concept. This combination provides for preferential traffic handling for higher-priority packets. It can selectively discard lower-priority traffic when the interface starts to get congested and provide differentiated performance characteristics for different classes of service.

## 3.7   Admission Control

To make appropriate internal and external admission control decisions and to configure edge devices correctly, each DiffServ domain is outfitted with a bandwidth broker (BB). When a sender signals its local bandwidth broker to initiate a connection, the user is authenticated and subjected to a local policy-based admissions control decision. On behalf of the sender, the bandwidth broker then initiates an end-to-end call setup along the chain of routers to be traversed by the flow. The bandwidth broker abstraction is critically important because it allows separately administered network clouds (possibly implemented with very different underlying

Table 3.6: Advantages of DiffServ over IntServ

|  | Integrated Services | Differentiated Services |
|---|---|---|
| Service differentiation | Individual flow | Aggregate of flows—Therefore simple and scalable |
| QoS approach | End-to-end | Per-hop—Therefore simple and coarse |
| Signalling protocol | RSVP | Bandwidth broker for absolute DiffServ |
| Scalability | Limited by no. of flows | Limited by no. of service classes |
| Service guarantees | Deterministic | Absolute or relative |

technologies and polices) to manage their network resources as they see fit. Within the QoS domain the bandwidth broker acts as a decision maker on admission control, service selection and service allocation. To do this the bandwidth broker should be able to estimate current load and performance information within its network.

The main functions of bandwidth broker are :

- collect and measure the local traffic load of its network
- estimate the demand level of each service for each IP bit pipe
- use all this information for admission control
- use bandwidth efficiently

## 3.8   Comparison of DiffServ and IntServ

DiffServ and IntServ are used to provide *Quality of Service* in IP networks. Both of these approaches have suitable applications. IntServ is more appropriate for relatively small or closed networks with high QoS requirements. DiffServ is better suited for a large network like the Internet. DiffServ is more scalable when it comes to fighting for resources. Therefore at present it is seen as the future technology. Advantages of DiffServ over IntServ can be seen from the Table 3.6.

The only advantage of IntServ over Diffserv is that IntServ is designed to guarantee end-to-end performance like a *circuit switched network*. This performance comes at the cost of implementing a reservation protocol and limiting scalability. DiffServ takes the less ambitious approach of providing QoS on a per-hop basis. The DiffServ model is scalable but does not provide the same level of performance guarantees.

# Traffic Management and QoS in Second and Third Generation Mobile Networks

## 4.1 Introduction

Existing wireless networks are mostly digital and support voice and data communication at a low bit rate of 9.6–16 kb/s. Fueled by the explosive growth of the Internet, applications are demanding that higher capacity, higher data rates, and advanced multimedia services be supported in the near future. The evolution to higher data rates and more advanced services occurs in two steps. The first step is the emergence of 2G + systems in which second-generation (2G) systems such as Global System for Mobile Communications (GSM) and IS-95 are extended to provide high-speed data communications either without changing the air interface or by using improved coding techniques. The second step is to provide higher capacity, data rates, and multimedia services. Wideband code-division multiple access (WCDMA) standard proposals such as the cdma2000 system include a greatly enhanced air interface to support wider bandwidths for improved capacity and higher data rates [Sar00]. UMTS is based on this new concept.

The first generation of wireless mobile communications was based on analog signaling. Analog systems, implemented in North America, were known as Analog Mobile Phone Systems (AMPS), while systems implemented in Europe and the rest of the world were typically identified as a variation of Total Access Communication Systems (TACS). Analog systems were primarily based on circuit-switched technology and designed for voice, not data. The second generation (2G) of the wireless mobile network was based on low-band digital data signaling. The most popular 2G wireless technology is known as Global Systems for Mobile Communications (GSM).

GSM systems, first implemented in 1991, are now operating in about 140 countries and territories around the world.

Today, GSM systems operate in the 900MHz and 1.8 GHz bands throughout the world with the exception of the Americas where they operate in the 1.9 GHz band. While GSM and other TDMA-based systems have become the dominant 2G wireless technologies, CDMA technology is recognized as providing clearer voice quality with less background noise, fewer dropped calls, enhanced security, greater reliability and greater network capacity. The Second Generation (2G) wireless networks mentioned above are also mostly based on circuit-switched technology. 2G wireless networks are digital.

2G wireless technology can handle some data capabilities such as fax and short message service at the data rate of up to 9.6 kbps, but it is not suitable for web browsing and multimedia applications. The virtual explosion of Internet usage has had a tremendous impact on the demand for advanced wireless data communication services. However, the effective data rate of 2G circuit-switched wireless systems is relatively slow -- too slow for today's Internet. As a result, GSM and other TDMA-based mobile system providers and carriers have developed 2G+ technology that is packet-based and increases the data communication speeds to as high as 384kbps. These are High Speed Circuit-Switched Data (HSCSD), General Packet Radio Service (GPRS) and Enhanced Data Rates for Global Evolution (EDGE) technologies.

## 4.2  General Packet Radio Service (GPRS)

The General Packet Radio Service (GPRS) is a new bearer service for GSM that greatly improves and simplifies wireless access to packet data networks, e.g., to the Internet. It applies a packet radio principle to transfer user data packets in an efficient way between mobile stations and external packet data networks. At present, data rates are too slow and the connection setup takes too long and is rather complicated. Moreover, the service is too expensive for users. From the technical point of view, the drawback results from the fact that current wireless data services are based on circuit switched radio transmission. At the air interface, a complete traffic channel is allocated for a single user for the entire call period. In case of bursty

Table 4.1: Coding Scheme (CS) parameters

| Coding scheme | Coding rate | Data rate[Kbps] |
| --- | --- | --- |
| CS-1 | $1/2$ | 9.05 |
| CS-2 | $\approx 2/3$ | 13.4 |
| CS-3 | $\approx 3/4$ | 15.6 |
| CS-4 | 1 | 21.4 |

traffic (e.g., Internet traffic), this results in a highly inefficient resource utilization. It is obvious that for bursty traffic, packet switched bearer services result in a much better utilization of the traffic channels. This is because a channel will only be allocated when needed and will be released immediately after the transmission of the packets. With this principle, multiple users can share one physical channel (statistical multiplexing). GPRS is developed to address these deficiencies. GPRS, or General Packet Radio Service, is considered to be one of the most important stages in making the Internet accessible via wireless telephones. GPRS is also sometimes know as 'mobile generation 2.5'. The European Telecommunications Standards Institute (ETSI) has approved four standards for GPRS: CS1 which transmits 9.05 kbps (kilobits per second) per timeslot, CS2 transmits 13.4 kbps, CS3 transmits 15.6 kbps and CS4 transmits 21.4 kbps [BVE99, Wap00] (Table 4.1).

One can expect a general model aimed at the general public, which will give four timeslots at 9.05 kbps (total 36.2 kbps). Other versions of CS2 and CS4 should also be available. The top-of-the-range, more expensive models offering faster transfer speeds would be aimed at providing professionals and the business community with the possibility of fast data transfer, where price is less of an issue. Packet switching means that GPRS radio resources are used only when users are actually sending or receiving data. Rather than dedicating a radio channel to a mobile data user for a fixed period of time, the available radio resource can be concurrently shared between several users. The actual number of users supported depends on the application being used and how much data is being transferred. Through multiplexing of several logical connections on one or more GSM physical channels, GPRS reaches a flexible use of channel capacity for applications with variable bit rate.

GPRS is extremely efficient in its use of scarce spectrum resources and enables GSM operators to introduce a wide range of value-added services for market differentiation. It is ideal for bursty type data applications such as e-mail, Internet access or WAP-based applications. GPRS is the packet mode extension to GSM. It uses the same air interface but with a new physical channel called a 52-multiframe, which is made of two 26 control multiframes of voice mode GSM. Packet mode control and data channels are mapped into different slots of the 52-multiframe, which takes 240 ms. 52-multiframe consists of 12 blocks (B0B11) of four frames to which several packet mode logical channels can be mapped, and four additional frames [Sar00]. As the transmission of long packets in packet-oriented systems reduces the economy of scale and increases the blocking probability, the work space of GPRS is primary:

- the frequent, regular transmission of short data packets up to 500 byte and
- the irregular transmission of short to middlesize data packets up to a few kbyte.

The basic approach to integrate the packet data service into the GSM standard represents the reservation and the logical subdivision of certain GSM channels. The number of channels allocated for GPRS is dynamically adapted to the workload situation in the respective cell. On the GSM broadcast channel BCCH the mobile subscriber is indicated on which frequency the GPRS is offered. That way up to eight traffic channels per frequency can be utilized by a subscriber.

### 4.2.1 Logical Architecture

In order to integrate GPRS into the existing GSM architecture, a new class of network nodes, called GPRS support nodes (GSN), has been introduced [BVE99]. GSNs are responsible for the delivery and routing of data packets between the mobile stations and the external packet data networks (PDN). Figure 4.1 illustrates the system architecture. A serving GPRS support node (SGSN) is responsible for the delivery of data packets from and to the mobile stations within its service area. Its tasks include packet routing and transfer, mobility management (attach/detach and location management), logical link management, and authentication and charging functions. The location register of the SGSN stores location information (e.g., current cell, current VLR) and user profiles (e.g., IMSI, address(es) used in the packet data network) of all GPRS users registered with this SGSN. A gateway GPRS sup-

Figure 4.1: Logical Architecture of GPRS [BVE99]

port node (GGSN) acts as an interface between the GPRS backbone network and the external packet data networks. It converts the GPRS packets coming from the SGSN into the appropriate packet data protocol (PDP) format (e.g., IP or X.25) and sends them out on the corresponding packet data network. In the other direction, PDP addresses of incoming data packets are converted to the GSM address of the destination user. The readdressed packets are sent to the responsible SGSN. For this purpose, the GGSN stores the current SGSN address of the user and his or her profile in its location register. The GGSN also performs authentication and charging functions. In general, there is a many-to-many relationship between the SGSNs and the GGSNs: A GGSN is the interface to external packet data networks for several SGSNs; an SGSN may route its packets over different GGSNs to reach different packet data networks.

There are two kinds of GPRS backbone networks:

- Intra-PLMN backbone networks connect GSNs of the same PLMN and are therefore private IP-based networks of the GPRS network provider.

- Inter-PLMN backbone networks connect GSNs of different PLMNs. A roaming agreement between two GPRS network providers is necessary to install such a backbone.

### 4.2.2 Parallel Use of Services

During a GPRS session circuit-switched services (speech as well as data) may still be initiated and used. Similarly, it is possible to send and receive GPRS data while carrying out a telephone call. Parallel use of these services is provided for *point-to-point* (PTP) as well as *point-to-multipoint* (PTM) services, causes. Contrary to the existing GSM the management of GPRS service profiles is based on the conception of the *multiple subscriber profile* (MSP). Thus, it is service specific, which means that a subscriber can activate every subscribed service separately.

### 4.2.3 QoS in GPRS Release 99

In GPRS Release 99 four different traffic classes are introduced, with different parameters specifying their QoS requirements (see Table 4.2). They are:

- *conversational*,
- *streaming*,
- *interactive*, and
- *background*.

Delay-sensitive services belonging to *conversational* class, e.g., do need absolute guarantees in terms of *guaranteed bitrate* and *transfer delay* attributes, while for *background* traffic none other than bit integrity is necessary.

## 4.3 Universal Mobile Telecommunication System (UMTS)

UMTS is one of the major new third generation mobile communications systems being developed within the framework which has been defined by the ITU and known as IMT-2000. UMTS will deliver pictures, graphics, video communications and other wide-band information as well as voice and data, direct to people who can

Table 4.2: End-user performance expectations for selected services belonging to different traffic classes [ETS00a]

| Traffic class | Medium | Application | Data rate | One-way delay |
|---|---|---|---|---|
| Conversational | Audio | Telephony | 4–25 kbit/s | < 150 ms |
| | Data | Telnet | < 8 kbit/s | < 250 ms |
| Streaming | Audio | Streaming audio (HQ) | 32–128 kbit/s | < 10 s |
| | Video | One-way | 32–384 kbit/s | < 10 s |
| | Data | FTP | - | < 10 s |
| Interactive | Audio | Voice messaging | 4–13 kbit/s | < 1 s |
| | Data | Web-browsing (HTML) | - | < 4 s/page |

be on the move. UMTS will build on and extend the capability of today's mobile technologies (like digital cellular and cordless) by providing increased capacity, data capability and a far greater range of services using an innovative radio access scheme known as WCDMA and an enhanced, evolving core network. It will provide data speeds of upto 2 Mbps, making portable videophones a reality. UMTS has the support of many major telecommunications operators and manufacturers because it represents a unique opportunity to create a mass market for highly personalised and user friendly mobile access to the Information Society.

**Spectrum for UMTS:** World Radio Communication Conference (WRC) 2000 identified the frequency bands 1885-2025 MHz and 2110-2200 MHz for future IMT-2000 systems, with the bands 1980-2010 MHz and 2170-2200 MHz intended for the satellite part of these future systems.

UMTS basic deployment is said to start in the year 2002. At present large amount of research is going on in this area. UMTS will be using the advantages of WCDMA over TDMA and FDMA. One of the main targets of UMTS is to provide mobile multimedia services with high bit rates. UMTS Terrestrial Air Interface(UTRA) is required to provide 2Mbps in limited environments, like in indoors and in small micro cells.

## 4.3.1 Requirements for QoS

**End User Requirements for QoS:** Generally, end users care only the issues that are visible to them. From the end-user point of view:

- only the QoS perceived by end-user matter
- the number of user defined/controlled parameters has to be as small as possible
- derivation/definition of QoS attributes from the application requirements has to be simple
- QoS attributes shall be able to support all applications that are used, a certain number of applications have the characteristic of asymmetric nature between two directions, uplink/downlink
- QoS definitions have to be future proof
- QoS has to be provided end-to-end

**Technical Requirements for QoS:**

- the UMTS QoS mechanisms shall provide a mapping between application requirements and UMTS services
- the UMTS QoS control mechanisms shall be able to efficiently interwork with current QoS schemes. Further, the QoS concept should be capable of providing different levels of QoS by using UMTS specific control mechanisms.
- multiple QoS streams per address should be possible
- the UMTS shall provide a finite set of QoS definitions
- the overhead and additional complexity caused by the QoS scheme should be kept reasonably low, as well as the amount of state information transmitted and stored in the network
- QoS shall support efficient resource utilization
- the QoS parameters are needed to support asymmetric bearers
- applications (or special software in UE or 3G gateway node) should be able to indicate QoS values for their data transmissions
- QoS behavior should be dynamic , i.e., it shall be possible to modify QoS parameters during an active session

- number of parameters should be kept reasonably low (increasing number of parameters, increase system complexity)

- user QoS requirements shall be satisfied by the system, including when change of SGSN within the Core Network occurs.

- QoS mechanism have to allow efficient use of radio capacity

- Allow independent evolution of Core and Access networks

- Allow evolution of UMTS network, (i.e., eliminate or minimize the impact of evolution of transport technologies in the wireline world).

### 4.3.2 QoS Architecture

To realize a certain network QoS a Bearer Service with clearly defined characteristics and functionality is to be set up from the source to the destination of a service. A UMTS bearer service layered architecture is depicted in Figure 4.2, each bearer service on a specific layer offers it's individual services using services provided by the layers below. utilize UMTS Bearer Service provides the UMTS QoS. The UMTS Bearer Service consists of two parts, the Radio Access Bearer Service and the Core Network Bearer Service. Both services reflects the optimized way to realize the UMTS Bearer Service over the respective cellular network topology taking into account such aspects as e.g. mobility and mobile subscriber profiles. The Radio Access Bearer Service provides confidential transport of signaling and user data between MT and CN Iu Edge Node with the QoS adequate to the negotiated UMTS Bearer Service or with the default QoS for signaling. This service is based on the characteristics of the radio interface and is maintained for a moving MT.

The Core Network Bearer Service of the UMTS core network connects the UMTS CN Iu Edge Node with the CN Gateway to the external network. The role of this service is to efficiently control and the backbone network in order to provide the contracted UMTS bearer service. The UMTS packet core network shall support different backbone bearer services for variety of QoS.

### 4.3.3 UMTS QoS Classes

When defining the UMTS QoS classes the restrictions and limitations of the air interface have to be taken into account. It is not reasonable to define complex
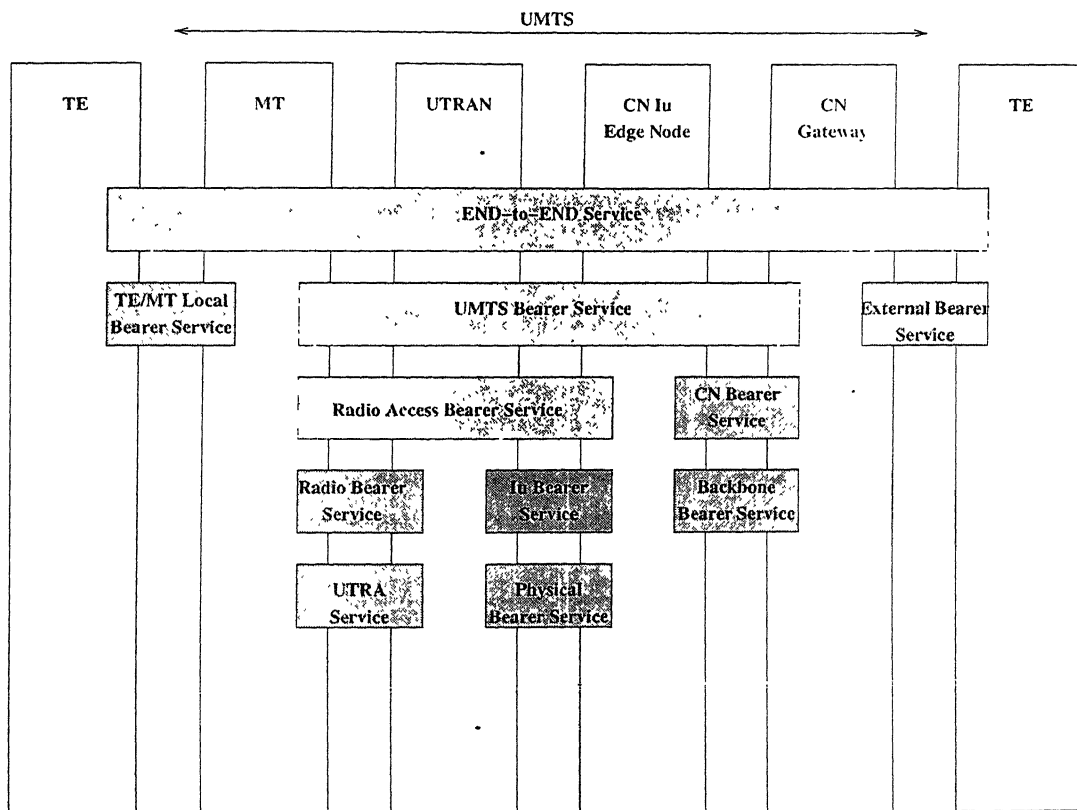
Figure 4.2: UMTS QoS Architecture [ETS99]

mechanisms as have been in fixed networks due to different error characteristics of the air interface. The QoS mechanisms provided in the cellular network have to be robust and capable of providing reasonable QoS resolution. In UMTS as in GPRS99 release we have four classes *conversational, streaming, interactive, background* [ETS99] (Table 4.3). The main distinguishing factor between these classes is how delay sensitive the traffic is: .Conversational class is meant for traffic which is very delay sensitive while Background class is the most delay insensitive traffic class. Conversational and Streaming classes are mainly intended to be used to carry real-time traffic flows. The main divider between them is how delay sensitive the traffic is. Conversational real-time services, like video telephony, are the most delay . sensitive applications and those .data streams should be carried in Conversational class. Interactive class and Background are mainly meant to be used by traditional Internet applications like WWW, Email, Telnet, FTP and News. Due to looser delay requirements, compare to conversational and streaming classes, both provide better error rate by means of channel coding and retransmission. The main differ-

Table 4.3: UMTS QoS Classes

| Traffic class | Fundamental characteristics | Example of application |
|---|---|---|
| Conversational class | -Preserve time relation (variation) between information entities of the stream<br>-Conversational pattern is stringent with low delay | voice |
| Streaming class | -Preserve time relation (variation) between information entities of the stream | streaming video |
| Interactive class | -Request response pattern<br>-Preserve payload content | web browsing |
| Background class | -Destination is not expecting the data within a certain time<br>-Preserve payload content | background download of emails |

ence between Interactive and Background class is that Interactive class is mainly used by interactive applications, e.g. interactive Email or interactive Web browsing, while Background class is meant for background traffic, e.g. background download of Emails or background file downloading. Responsiveness of the interactive applications is ensured by separating interactive and background applications. Traffic in the Interactive class has higher priority in scheduling than Background class traffic, so background applications use transmission resources only when interactive applications do not need them. This is very important in wireless environment where the bandwidth is low compared to fixed networks.

### 4.3.4   UMTS-Internet Interworking

The main goal of UMTS QoS mechanisms is to create a future proof concept that will provide means to transport different types of data with different QoS requirements. Thus the interworking of UMTS and existing/evolving network technologies has to be ensured. In the case of Internet applications, the selection of the class and

| UMTS QoS class | DiffServ class | Reason |
|---|---|---|
| Conversational class | Expedited Forwarding class | As it requires low latency and jitter and EF class guarantees a minimun service level |
| Streaming class | Assured Forwarding class 4 | As it requires low variation of delay i.e. stringent jitter requirements |
| Interactive class | Assured Forwarding class 3 | As it requires low latency (but not as low as in conversational class) |
| Background class | Assured Forwarding class 2 or class 1 or best effort (class 0) | As there is no specific requrement for this class except reliability, it can be given to AF class2 or AF class 1 or even to best effort service |

Figure 4.3: A proposal for mapping UMTS classes with DiffServ classes

appropriate traffic attribute values is made according to the Internet QoS parameters. Internet applications do not directly use the services of UMTS but they use Internet QoS definitions and attributes, which are mapped to UMTS QoS attributes at the interface [ETS99]. Currently there are three main Internet QoS concepts, namely Integrated Services and Differentiated Services and MPLS. In this thesis an IP Core Network Model is set up which is based on Differentiated Services which is then connected to the GPRS Radio Network. DiffServ requires that there is either one QoS profile for each traffic type or alternatively the priority and traffic type information is included in the data packets. A proposal for mapping UMTS classes with DiffServ classes is shown in Figure 4.3.

## 4.3.5  Third Generation Core networks

In order to provide the bandwidth flexibility depending upon end-user requirements, 3G networks will use a packet-based approach, rather than the circuit-switched technology of second generation systems. The network architecture due for march 2001 release of 3GPP standardization consists of three distinct layers:

**Application layer:**  It is where the end-user applications reside. In modern networks, applications are implemented in the mobile terminals and on dedicated application servers in the network.

**Network control layer:**  It houses a number of network servers of different types. These servers are responsible for controlling such aspects as mobility management, set-up and release of calls and sessions requested by the end-users. Only control communications are handled here and no user plane streams pass through.

**Connectivity layer:**  The connectivity layer is a pure transport mechanism that is capable of transporting any type of information like voice, data and multimedia streams. Its backbone architecture incorporates core and edge equipment. The core equipment transports aggregated traffic streams between the different nodes at the edges of the backbone. As a rule, core equipment is a backbone router or backbone switch that handles traffic streams either according to very simple classification principles, or according to routes that the network operator has predefined by means of traffic engineering. Edge equipment collects customer specific data and statistics for accounting and billing purposes, and provides the single bit pipes that guarantee an appropriate quality of service. The edge equipment is usually a media gateway, which operates under the full control of the nodes in the network control layer. Connectivity layer solutions can be based on either asynchronous transfer mode (ATM) transmission or Internet protocol (IP) transmission. However, in the long run, the role of ATM as a transport technology is expected to diminish, leaving IP transmission as the predominant technology. This Internet protocol will have added *quality of service* features of *DiffServ and MPLS* [Wit00, FHPW00].

# The GPRS and EGPRS Simulator GPRSim

## 5.1 Simulation Environment

In this chapter we will talk about the *General Packet Radio Service simulator* that is developed at the chair of communication networks, rwth-aachen. This simulator is made of many modules, some of these modules are discussed in this chapter with their functionality. The *General Packet Radio Service* (GPRS) and *Enhanced General Packet Radio Service* (EGPRS) simulator GPRSIM has been developed as a pure software solution utilizing the programming language C++. Up to now, models of mobile station, base station, and SGSN are implemented in the simulation environment. Interfaces are offered for the simulator to be upgraded by additional modules.

For implementation of the simulation model in C++ the *Communication Networks Class Licrary*(CNCL), which was developed at the Chair of Communication Networks (COMNETS), is used. This class library allows an object-oriented structure of programs and is especially applicable for event-driven simulations. The complex protocols like LLC, RLC/MAC and the Internet load generators are specified with the SDL, are translated to C++ by the code generator SDL2CNCL [SOL97], and are finally integrated into the simulator.

Different from usual approaches to building a simulator, where abstractions of functions and protocols are being implemented, the approach of the GPRSIM is based on detailed implementation of the standardized protocols. This enables studying the behaviour of GPRS and EGPRS in a realistic way.
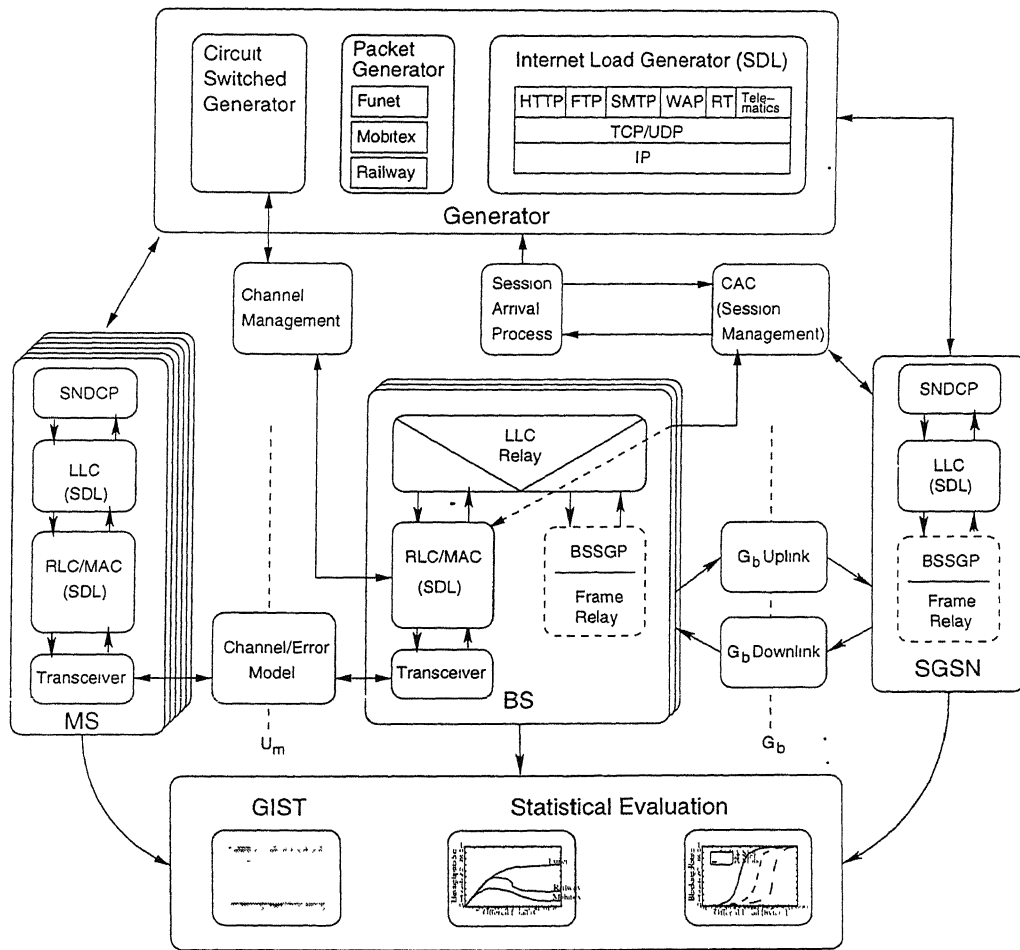
Figure 5.1: The structure of the (E)GPRS simulator GPRSIM

## 5.2 Structure of the GPRS Simulator

The software architecture of the GPRSIM and the information flow between the modules is shown in Figure 5.1. The figure represents the state of the simulator before modification carried out in this thesis. Modified software architecture of GPRSim is shown in Figure 6.2 in the next chapter.

The GPRSIM comprises the modules MS, BS, SGSN, the transmission links, the load generators, session control modules, a Web interface, and a module for statistical evaluation.

The MS, BS, and SGSN modules contain the implementations of the respective protocol stacks. Transmission links are represented by simple error models. While the $G_b$ interface is regarded as ideal, block errors on the radio interface $U_m$ can be

simulated based on look-up tables, which map a C/I value to a BLEP (block error probability).

The generator comprises Internet applications like *www, smtp. ftp* and a *Circuit Switched* (CS) generator which models voice traffic coexisting to the packet-oriented data traffic.

The module *Channel Management* supervises the physical GSM channels available in the cell and allocates channels for the GPRS Radio Resourse Management entity. either as fixed packet data channel (PDCH) or on demand.

The output of the simulator comprises a graphical presentation of the protocol cycle and statistical evaluation of performance measurements. Some modules are presented with their functionality and interactions.

## 5.2.1   Load Generators

The GPRSIM comprises an Internet load generator based on Internet traffic models for the applications HTTP, FTP, and SMTP. The parameters of these models coming from American traffic measurements are updated by parameters given by ETSI UMTS prognoses for the behaviour of mobile Internet users [ETS97]. Traffic models for WAP have been developed and those for *Telematics* and *Real Time* applications are under development.

Besides the Internet load generator a Circuit Switched generator to model coexisting voice traffic is also implemented. All generators communicate with the other GPRSIM modules via the *SDL process environment*.

### 5.2.1.1   Internet Load Generator

The conception of the Internet load generator is the specifications of a system *Client* that behaves like a user running an Internet session, and a system *Server* representing the peer entity for the Internet service [Stu99]. Underneath the generator the protocols TCP and IP are included. IP packets are sent to and received from the environment by server and client, respectively.

Following, a description of the applications WWW, E-mail, and FTP, and their handling utilizing mathematical descriptions for generating protocol-specific traffic, will be provided. The parameters of these models are derived from American traffic
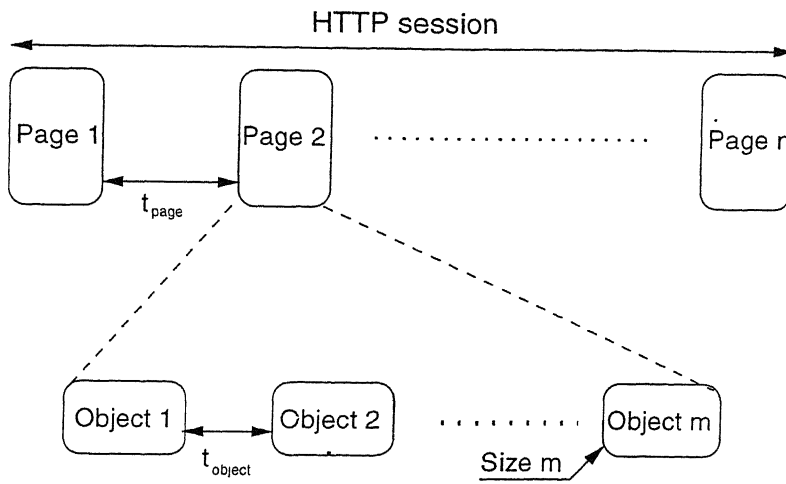
Figure 5.2: Sequence chart and important parameters for an HTTP session

measurements are updated by parameters given by ETSI UMTS prognoses for the behaviour of mobile Internet users [ETS97].

**The HTTP model:** To describe this load model's parameters, Figure 5.2 depicts the sequence chart for a complete session on a WWW browser.

When describing the mathematical form of the load generation, the value $u$ between 0 and 1 will always reflect a random number generated by an equal-distributed random number generator. From this value the others are calculated respective the rules below.

HTTP sessions consist of requests for a number of *pages*. This value follows a *geometric distribution* with a mean value of $\overline{n} = 5$. In mathematical nomenclature a generation of geometric-distributed values is designated by

$$ n = \frac{\ln(u)}{\ln(1-p)} \tag{5.1} $$

Herein, the parameter $p$ reflects the *probability of success* and is determined from the mean value $\overline{n}$

$$ p = \frac{1}{\overline{n}} \tag{5.2} $$

But not only the number of pages describes the behaviour, the delay between two pages must also be defined. In the following definition the time $t_{page}$ between two pages is defined as the delay between complete reception of the previous page until start of transmission of the new one; refer to Figure 5.2. This value heavily depends on the user's behaviour while browsing Web con-

Table 5.1: Model parameters of an HTTP session

| Parameter | Distribution | Mean |
|---|---|---|
| Pages per session | geometric | 5.0 |
| Intervals between pages [s] | negative-exponential | 12.0 |
| Objects per page | geometric | 2.5 |
| Object size [byte] | $\log_2$-Erlang-k | 3700 |

tent. The measurements in [AW95] result in a description for $t_{page}$ according to a *negative-exponential distribution* with a mean value of $\overline{t_{page}} = 33\,\text{s}$. ETSI, on the other hand, proposes $\overline{t_{page}} = 12\,\text{s}$ [ETS97].

Therefore. the generation of random numbers with negative-exponential behaviour follows

$$t_{page} = -a \cdot \ln(u) \tag{5.3}$$

with

$$a = \overline{t_{page}} \tag{5.4}$$

Referring to Figure 5.2, further differentiation has to be done. A page can be divided into a couple of *objects*. Three parameters need to be modeled to generate these objects. First of. all, the number $m$ of generated objects has to be specified. A mean value of $\overline{m} = 2.5$ objects per page following a geometric distribution is taken. Consequently, this parameter follows equation 5.1.

The delay between two objects is given by a constant value of $\overline{t_{object}} = 0\,\text{s}$.

Now, the *size* of each object has to be defined using an approximation. The amount of data for each object transferred is in reality not only defined by the object size but also by additional header information. This additional data will be neglected since it represents only a minor fragment of the total amount. [AW95] proposes a *$\log_2$-Erlang-k distribution* with a mean of $\overline{m} = 3700\,\text{byte}$. For the generation of object sizes the value $k$ is needed and determined to $k = 24$. The mathematical representation is given to

$$m = \frac{\overline{m}}{k} \cdot \ln\left(\prod_{i=0}^{k} u_i\right) \tag{5.5}$$

Table 5.1 gives an overview of the parameters describing the HTTP model.

Table 5.2: Model parameters of an FTP session

| Parameter | Distribution | Mean | Variance |
|---|---|---|---|
| Total bulk of data [byte] | $\log_2$-normal | 32768 | 10000 |
| Bulk of data per object [byte] | $\log_2$-normal | 3000 | 1000 |
| Interval between connections [s] | $\log_{10}$-normal | 4 | 2.55 |

**The FTP model:** The model introduced in this section represents a unidirectional data source. It describes an object transmission from an FTP server to an FTP client. The examinations are based on extensive WAN measurements documented in [Pax94]. Over 95 % of the measured FTP data connections are performed by a GET command. Therefore, it is sufficient to regard data transfer from server to client. Traffic offer originated by FTP control connections is not considered in the following model. Parameters describing an FTP session are:

- the total bulk of data per session
- the size of each transferred object (FTP data connection)
- the interval between two object transmissions

The bulk of data per session characterizes the duration of a session. In [Pax94] a *$\log_2$-normal distribution* is presented for the total bulk of data per session. As parameters the mean value of $\overline{x}_{arith.} = 32768$ byte and the standard deviation of $\sigma_{x,arith.} = 10000$ byte are suggested. For the bulk of data per object – per FTP data connection – a *$\log_2$-normal distribution* is detected, likewise. This distribution is characterized by the mean of $\overline{x}_{arith.} = 3000$ byte and a standard deviation of $\sigma_{x,arith.} = 1000$ byte. The given number of objects includes both the file transfer itself and, e.g., the listing of a directory. To describe the interval between two object transmissions a model created in [PF95] can be used. Intervals between two objects should follow a *$\log_{10}$-normal distribution*. Mean of the measurements is $\overline{x} = 4$ s and the variance is $\sigma_z^2 = 2.55$ s. These parameters are summarized in Table 5.2.

**The SMTP model:** The following SMTP model describes the load arising with the transfer of a message downloaded from a mail server by an electronic mail user. The bidirectional phase at the beginning of an SMTP session is not recreated

Table 5.3: Model parameters of an SMTP session

| Bulk of data [byte] | Distribution | Mean | Variance |
|---|---|---|---|
| E-mail size | $\log_2$-normal | 10000 | 1000 |
| Base quota | constant | 300 | — |

by this model. It only describes the bidirectional phase of the message transfer from SMTP server to SMTP client itself.

The only parameter is the bulk of data per SMTP session (E-mail). Measurements from [Pax94] are characterized by *$\log_2$-normal distribution* and a fixed added quota of 300 byte. The parameters of this distribution are shown in Table 5.3.

# Implementation of IP core network model

In this thesis an *IP core Network Model* is developed in *Specification and Description Language* (SDL) and then interconnected with the *General Packet Radio Service simulator* (GPRSim). An IP core Network model based on DiffServ consists of edge routers and interior routers. Both these routers are developed as *SDL modules*. and then are converted into C++ class by means of *SDL2CNCL* code generator developed at *Institute of Communication Networks, RWTH-Aachen* where CNCL is Communication Networks Class Library. All SDL modules communicate with each other and other modules via *SDL Process Environment.*

## 6.1   Edge Router Model in SDL

Figure 6.1 shows the Edge Router Model in SDL. Its main work is to carry out Connection Admission Control (CAC) and to classify and condition the incoming packets.

**Block Cl-Co-Dr:** This SDL block is nothing but a traffic classifier and traffic conditioning block. All the incoming packets are classified and are then given a particular DiffServ field codepoint depending on their priority. It converts the incoming data packets into IP packets. All the default data packets (packets for which connection admission control is not asked for) are given DiffServ priority level zero i.e. the best effort class.

**Block CAC:** All the connection admission control requests coming from the GPRS Radio Network arrive at connection admission control (CAC) block. It then forwards the CAC requests to traffic classifier and conditioner block (Cl-Co-Dr) for necessary action. This block then gives it a DiffServ field code point and a DiffServ priority level of 7 (highest priority class). This request is then converted into an IP packet and send looking for resources that it will need if
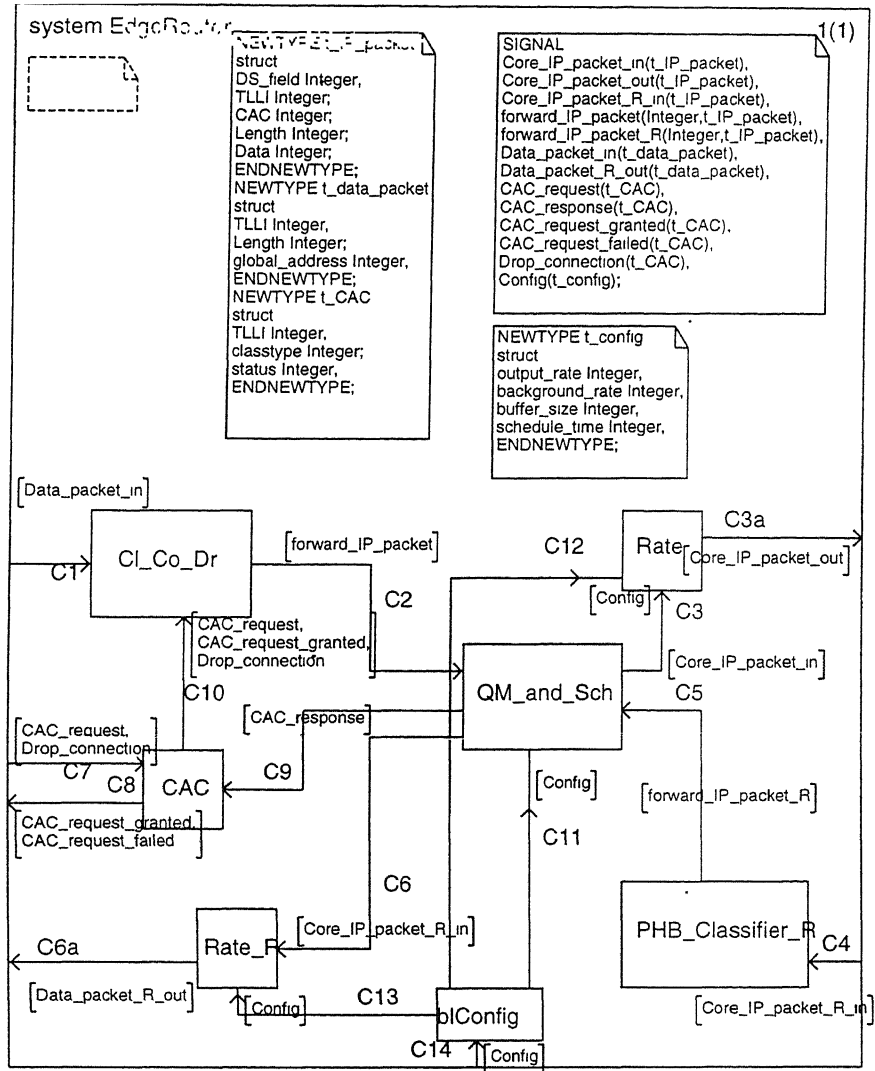
Figure 6.1: The Edge Router Model in SDL

the connection is granted. As this is of high priority it will be routed faster
than any other class. If the resources are available and the request is granted
then CAC sends back *CAC granted* signal back towards GPRS Radio Network.
If the request for particular priority fails then CAC checks for the next priority
level.

**Block blConfig:** It is used to configure the routers before the start of the simula-
tion. Default configuration values are shown in Table 6.1.

Table 6.1: Default Configuration Values

| Parameter | Value |
|---|---|
| output link rate | 2Mbps |
| background rate | 1Mbps |
| buffer size for each priority level | 0.1 MB |
| schedule time | 20ms |

**Block PHB-Classifier-R:** This block is per-hop behavior classifier block. It is used to classify packets only on the basis of DiffServ field as per the given DiffServ standard.

All the interior routers of the model contains only queue management and scheduling block and per-hop behavior classifier block. As in DiffServ it is said that the interior routers should be asked to perform only simple functions. Complexity should be added only on the edges of the network and the interior should be free from it. This is also because in an end-to-end connection there will be many interior routers but only few edge (boundary) routers.

## 6.2 GPRS simulator GPRSim with IP core network model

The General Packet Radio Service (GPRS) Core Network (see Figure 6.2) is developed as a part of this thesis and then it is interconnected to the already existing GPRS Radio Network. Thus now a connection is granted by looking at resources in both Radio as well as Core Network. Once a connection is granted packets are transported from host to destination and vice versa. In uplink the packets from host (Mobile Station) first go to the Radio Network and then to the Core Network, from where it goes to its destination. Similarly in downlink data packets from destination reach the host (Mobile Station) via Core and Radio Network respectively.

In this thesis IP routers are developed where edge routers are given to Serving GPRS Support Nodes (SGSN) and Gateway GPRS Support Node (GGSN). SGSN and GGSN are connected to each other via a number of Interior routers as shown in Figure 6.3 and Figure 6.4.
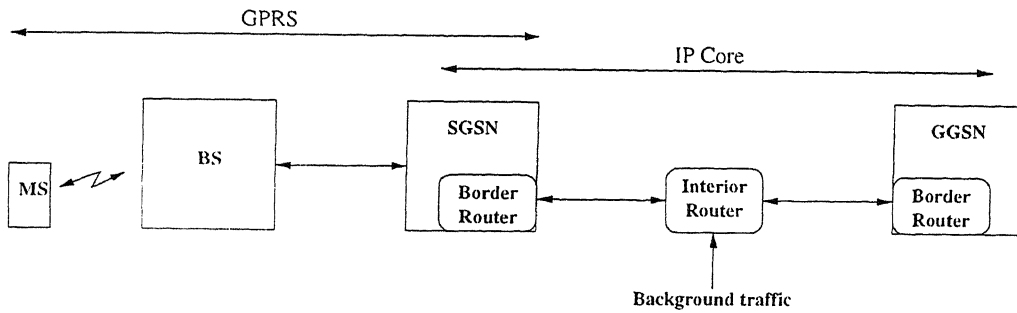
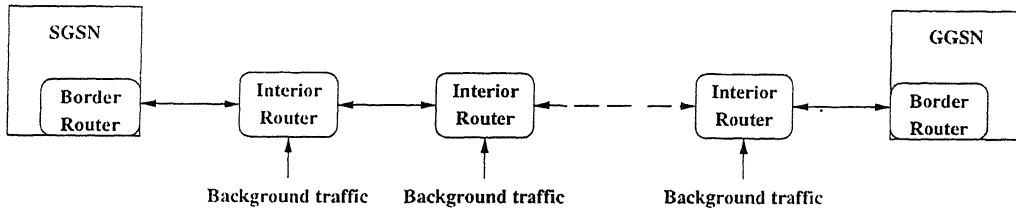Figure 6.3: IP core model with one Interior Router



Figure 6.4: IP core Model with a cascade of Interior Routers

## 6.3 Downlink data flow in the simulation model

In the model all SDL modules communicate with each other with the help of SDL Process Environment (see Figure 6.5). Whenever SDL Server (destination) has a it data packet to be sent to SDL Client (host or Mobile Station) it sends it to SDL Process Environment. SDL Process environment sends it to the Generator from where it is sent to the Core Network Interface. Core Network Interface sends it back to SDL Process Environment so that it can be sent to SDL system Edge Router GGSN. Thus from here it is sent to SDL system Edge Router GGSN (i.e. to the border router of the core network nearest to the destination). From GGSN the *data packet* is sent to SDL system Interior Router (next in line router after border router) via SDL Process Environment.

Then the *data packet* moves from one interior router on its path to next via the SDL Process Environment towards the interior router nearest to the SGSN. From here it is sent to SDL system Edge Router SGSN (i.e. other end border router of the Core Network). Once *data packet* comes out of the SDL system Edge Router SGSN it is out of the Core Network. Then it goes to the Core Network Interface via the SDL Process Environment. From Core Network Interface it is given to the Generator. Generator then sends the *data packet* towards the client via the GPRSim
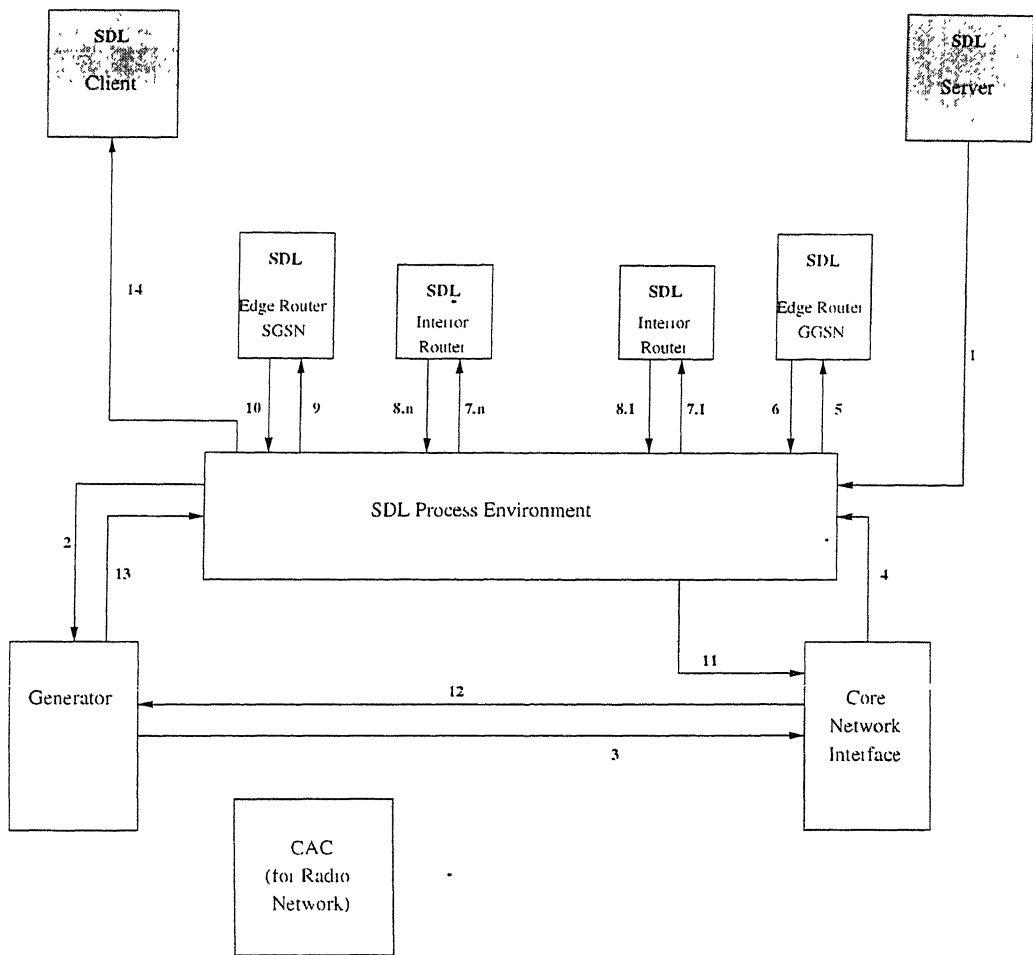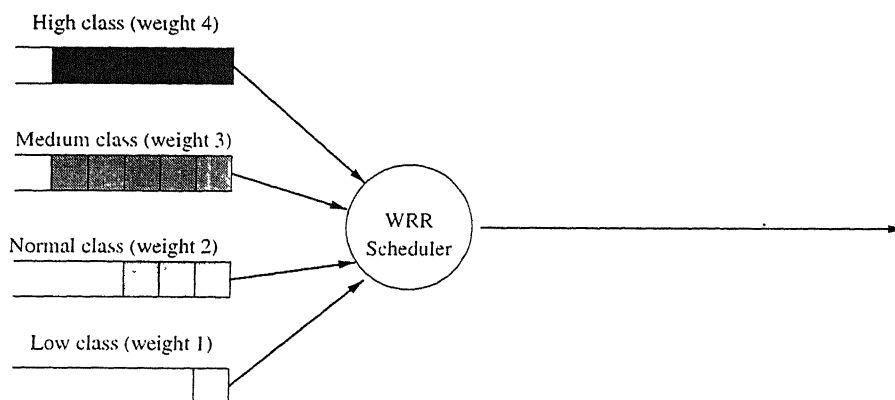
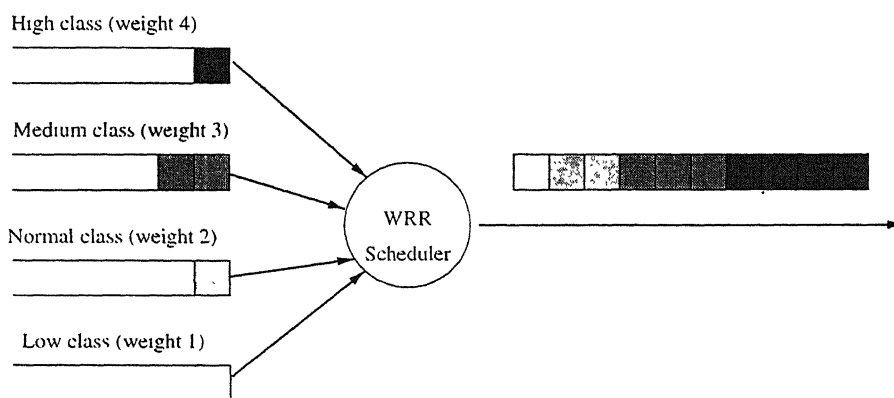Figure 6.5: Downlink data flow in the implementation

Radio Network taking help from the SDL Process Environment. This way downlink data flows in the *General Packet Radio Service simulator* (GPRSim) with added features of Core Network.

## 6.4 Weighted Round Robin Scheduling

In both edge routers and interior routers the scheduling algorithm that is used is weighted round robin (WRR) scheduling. The reason for choosing it was the DiffServ standard which says that the interior routers should be free from complex methods and work performed by interior routers should be a minimum. This is because in an end-to-end connection there will be a number of interior routers involved. Therefore its architecture should be simple.
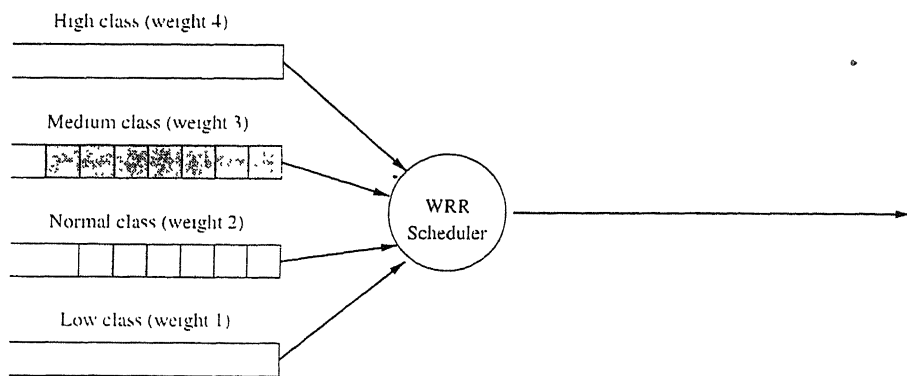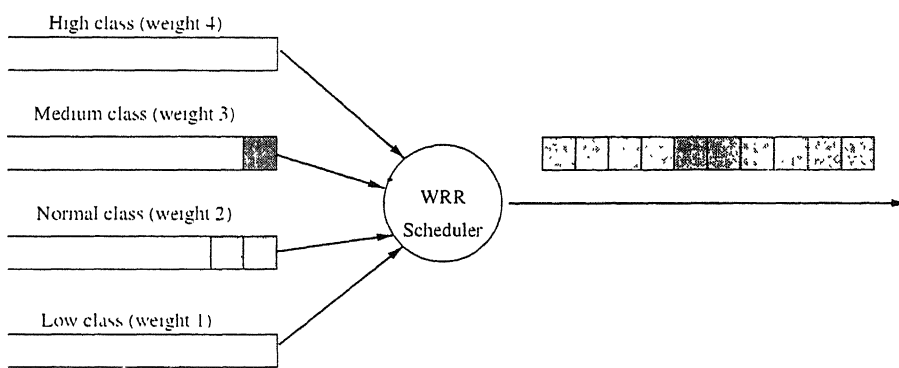
Figure 6.6: Weighted Round Robin example when no queue in empty

Another method of scheduling that can be used was weighted fair queuing (WFQ) but as it would require that each interior router should also know about the flow states of the packets passing through it, it was not implemented. By knowing about flow state we mean that each interior router should know the *number of flows* of each priority class passing through it at every schedule time instant. In the model the interior routers do not care about the number of flows of each priority class passing through it. Thus they are not directly concerned with the connection admission control. Whereas in Resource Reservation Protocol (RSVP) all routers on the path must know about the end-to-end connections they service and how much bandwidth each connection requires.

(a)



(b)

Figure 6.7: Weighted Round Robin example when some queues are empty

Weighted round robin (WRR) scheduler schedules packets depending upon the weights of each priority class. In the model DiffServ best effort class (class 0) is given the weight 1, assured forwarding (AF) class 1 a weight 2, assured forwarding (AF) class 2 a weight 3, assured forwarding (AF) class 3 a weight 4, assured forwarding (AF) class 4 a weight 5, expedited forwarding (EF) class a weight 6, DiffServ class 6 a weight 7 and DiffServ class 7 a weight 8. Between schedule time the bandwidth is divided into these DiffServ classes depending upon their weight. It serves each queue in a round robin manner, and for each turn, a number of bits corresponding to the queue's weight is pulled out from the queue. If one or more classes is not using its full allocation, then the unused capacity is distributed to the other classes according to their weighting. The bandwidth is distributed among only those classes that do not have their buffers empty at the starting of every schedule time (see Figure 6.6

and Figure 6.7). Also even after allocation of bandwidth if the higher class is not using its quota then its quota is given to the next lower class. What happens is that at every schedule time instant the bandwidth is divided based on number of priority classes present at that moment. It doesn't take into account what has happened in the previous schedule time instant. Thus based on this methodology the output link is given data packets till the arrival of next schedule time instant.

## 6.5   Admission Control

In the IP core network model implemented the edge routers are given the function of performing admission control, manage network resources and configure the network. Thus in a sense bandwidth broker which is responsible for connection admission control in DiffServ network architecture resides in the edge router. Whenever a request arrives at the edge router from the mobile radio access network, he checks whether he is capable of taking the connection or not depending upon the networks resources available at that instant. This he does by sending a request to all the core routers and waiting for a response from them. Once he gets their response he makes his decision to grant or drop the request. But unlike RSVP this method of admission control does not allocate bandwidth throughout the path from source to destination. Thus here the core routers does not have any idea about connection admission and therefore they do not have to maintain a per flow state of each admitted flow.

# Simulation Results

The enhanced GPRSIM to which the IP core network model module is added has to be validated by means of suitable simulations. This chapter comprises an estimation of the gain in capacity and performance achievable through appliance of *Differenti-ated Services* in the IP core network model developed in this thesis as compared to best effort Internet model for the mobile core network.

## 7.1 Simulation Scenarios

In the initialization phase input parameters for the simulations are read from a parameter file. These values characterize the simulation scenarios.

### 7.1.1 Cell Configuration

As GPRS applications will be the minority compared with circuit switched applications within the GSM in the near future, only cell configurations up to eight packet data channel (PDCH) on demand are regarded.

### 7.1.2 Traffic Generation

For network capacity planning the system performance during the busy hour has to be regarded. Here the input load is characterized by the number of Internet sessions running. The system performance measurements in the framework of this thesis can be compared with the number of sessions supposed for a busy hour situation. The GPRSim produces plausible results for load situations up to few tens of sessions running simultaneously. So here all results are limited to twenty mobile stations running Internet sessions.

The traffic generated by Internet applications is recreated by the Internet Load Generator described in chapter 5. The read time per WWW page can be taken from 12 s to 30 s, 60 s and 120 s.

### 7.1.3   IP core network model

In the IP core network model which is based on differentiated services the output link rate i.e. the rate at which the data is transfered from one router to another. the background traffic which is given to the best effort DiffServ class can be configured. The default values are 2Mbps for output link rate and 1Mbps for background rate. The router buffer size of each priority level and the schedule time can also be configured.

## 7.2   Simulation Parameters

The scenarios set up for the following simulation series first of all differ in the number of *Transceiver* (TRX) units in a cell. One TRX unit corresponds. depending on the scenario, to 6–8 physical GSM channels that have to be shared by *packet switched* and *circuit switched* traffic. There is a maximum of eight packet data channels (PDCH) available for packet switched traffic within each scenario. Thus, the number of on-demand packet data channels (PDCH) is calculated from: 8 − number of fixed PDCH. An exception is given by the scenario with 1 TRX; in this case, there are only 6 channels in the cell. The number of PDCH allocated on-demand is then: 6 − number of fixed PDCH.

The maximum possible number of on-demand packet data channels (PDCH) is made available for GPRS, as long as they are not required for the transmission of circuit switched traffic. When a request for a circuit switched transmission resource cannot be served by a free traffic channel (TCH), there is instantly detracted a channel from the pool of on-demand PDCH and utilized as a traffic channel.

The *packet switched* traffic to be generated for GPRS is in the ratio of 3:7 for *www* and *e-mail* sessions respectively. For the creation of coexisting *circuit switched* traffic load the CS generator introduced in chapter 5 is utilized. The relative quantity of CS traffic is the same throughout all simulations.

The evaluation of system capacity and QoS is then done on basis of following parameters.

**Mean IP throughput per user:** For this, the IP packet throughput is measured on a per-train basis (for both uplink and downlink), i. e., IP packets that belong to the same request and are therefore transmitted without interruption, e. g., the data of a single object on an HTML page. This is a crucial parameter for the *quality of service* of the system from a user's point of view. The statistical evaluation of this measure is done by counting the amount of IP bytes received in each TDMA frame period if a packet train is running. Each value divided by the TDMA frame duration represents a value in the evaluation sequence and is written to the evaluation container class. At the end of the simulation, the mean throughput is calculated from this evaluation sequence.

**Mean IP throughput per cell:** The mean IP throughput in the cell is calculated from the total IP data transmitted on all channels and for all users, divided by the simulation duration. In the loss-free system regarded, this will be equal to the offered traffic per cell.

**Mean IP packet delay:** The end-to-end delay of IP packets is evaluated by means of time stamps the packets are given when they move out from sender towards the receiver. When the packet arrives at the receiver, the difference of the current time and the timestamp value is calculated. This value is one entry of the evaluation sequence which consists of all such entries.

**IP packet delay (95-percentile):** 95-percentile means that 95 % of all IP packets arrive at the receiver within this time span. Only 5 % of the packets are delayed longer during transmission.

**Mean number of PDCH available:** As long as there is none of the PDCH allocated on-demand needed to carry circuit switched traffic, the maximum number of PDCH provided for the respective scenario is available for the transmission of packet switched data, i. e., eight for more than one TRX. The number of channels available for GPRS averaged on the total simulation duration is taken for evaluation.

**Number of sessions per hour:** The number of sessions within an hour delivers a clue for the statistical reliability of the simulation. High-load simulations

of large number of MS can only be run for approximately a half of the run-time of the remaining simulations due to extraordinary run-time requirements. However, they possess higher values for the number of sessions generated within this time and may thus be regarded statistically reliable.

## 7.3 Simulation Results with Differentiated Services

Several graphs are plotted in next few pages and their simulation scenario can be seen from Table 7.1. In the General Packet Service simulator (GPRSim) traffic models for *interactive* (http) and the *background* (email) traffic exists. Thus the area for measuring performance was small but the results that came out shows clear sign of great advantage of differentiated services architecture for the 3G mobile core networks. Even though what was being observed was end-to-end performance of only two classes i.e. *interactive* and *background*, the results clearly satisfied the expectations. It is believed that the results will hold true in the presence of other traffic classes *(conversational, streaming)* as well.

Figure 7.1 shows the relationship between downlink IP throughput per user in Kbps with respect to number of mobile stations. As can be seen from the graph the downlink IP throughput per user for interactive traffic (http) is far better than the background traffic (email). This distinction becomes even clear when number of mobile stations are increased or in other words when the traffic increases. By carefully aggregating a multitude of QoS-enabled flows into a small number of aggregates that are given a small number of differentiated treatments within the network, DiffServ eliminates the need to recognize and store information about each individual flow in core routers. This basic trick to scalability succeeds by combining a small number of simple packet treatments with a larger number of per-flow policing policies to provide a broad and flexible range of services. Similarly from Figure 7.2 we can again see the same advantage of interactive traffic over background traffic in the case of uplink IP throughput per user.

Figure 7.3 is the plotted graph between downlink IP packet delay in ms and the number of mobile stations. It is seen that the delay for high priority interactive class is far less as compared to the low priority background class. This becomes more dominant as the number of mobile stations are increased. Thus as congestion

Table 7.1: Parameters for simulations with differentiated services in the IP core network of GPRSim

| Parameter | Value |
| --- | --- |
| Number of transceiver (TRX) | 3 |
| Model time | One hour |
| Number of MS | 1–20 |
| Packet data channels (PDCHs) fixed | 0 |
| Traffic channels (TCHs) | 21 |
| Packet data channels (PDCHs) on demand | 8 |
| Offered CS traffic [Erlang] | 12.838 |
| $C/I$ ratio [dB] | 12 |
| Session inter-arrival time [s] | 12 |
| Relation $www/email$ (Interactive/background sessions) | 3/7 |
| HTTP read time [s] | 12 |
| Pages per session | 5 |
| Objects per page | 2.5 |
| HTTP object size [byte] | 3700 |
| E-mail size [byte] | 10000 |
| Number of standard subscribers [%] | 100 |
| IP core output rate | 2Mbps |
| IP core background traffic rate | 1Mbps |
| Buffer size per priority class | 0.1MB |
| schedule time | 40 ms |

increases, traffic in a given class will experience performance degradation due to congestion later than traffic in a lower class. Additionally, high precedence traffic may be guaranteed to experience less queuing delay than low precedence traffic. Users may still contract for a specified profile at a specified precedence level, but there is no way to characterize the service that a flow receives in an absolute sense. Same can also be seen from Figure 7.4 for the delay in the uplink direction. It is therefore argued that DiffServ appears to be remarkably good fit for mobile core networks as it will be dealing with large number of flows at every given instant. Mobile core networks should therefore explore DiffServ in wide area testbed.
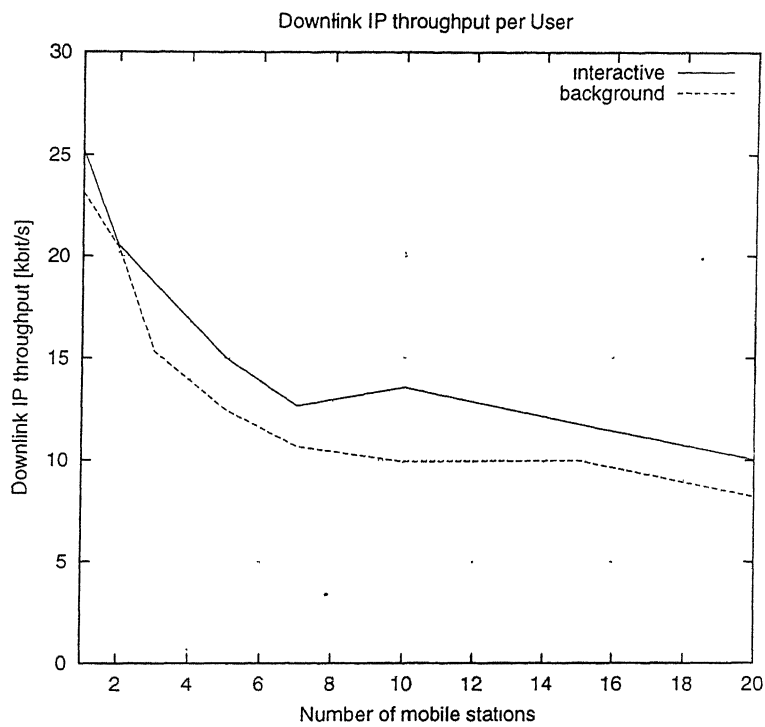
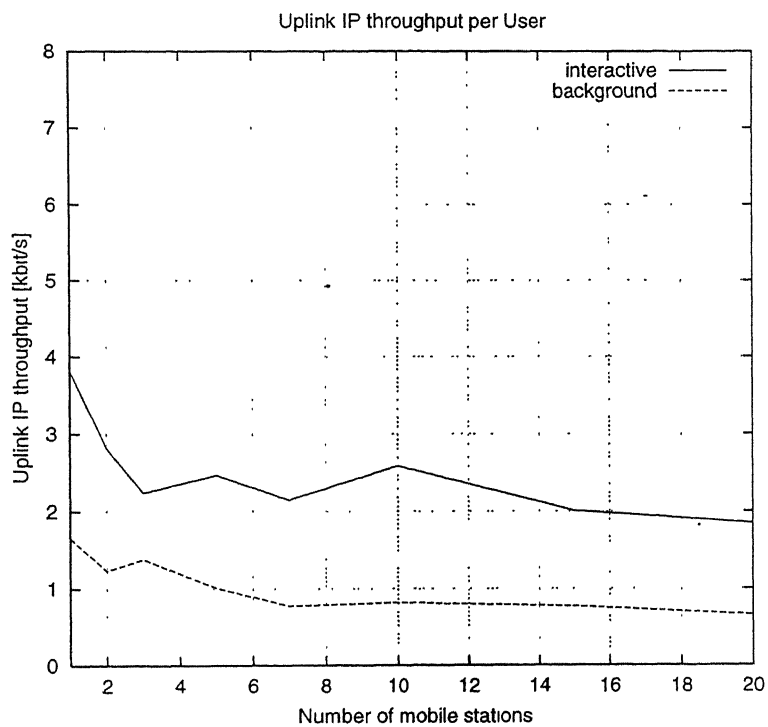Figure 7.1: Downlink IP throughput per user
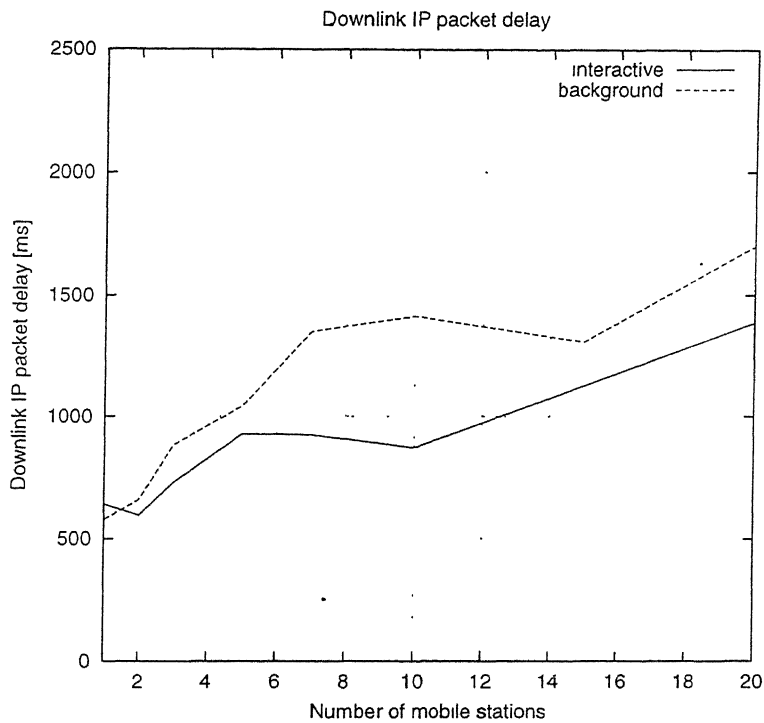


Figure 7.2: Uplink IP throughput per user

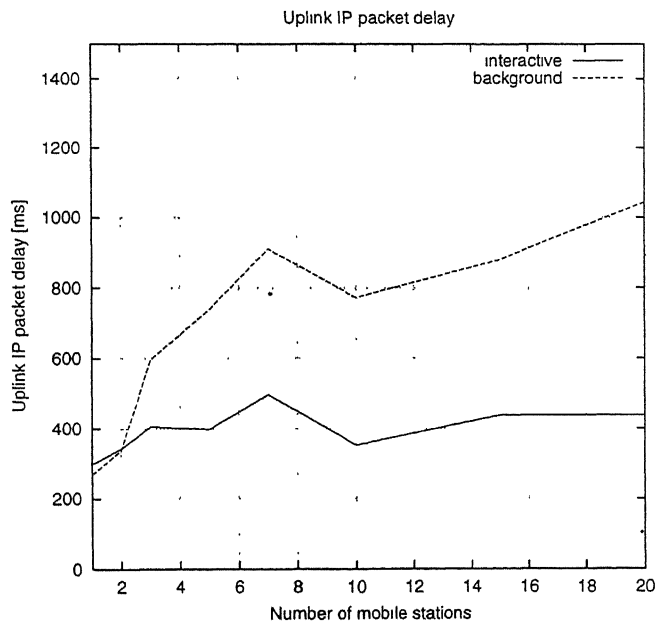Figure 7.3: Downlink IP packet delay



Figure 7.4: Uplink IP packet delay

Similarly from Figure 7.5 it can be inferred that no matter what the load situations of the mobile core network is, the high priority class will always have better performance level as compared to lower priority level.
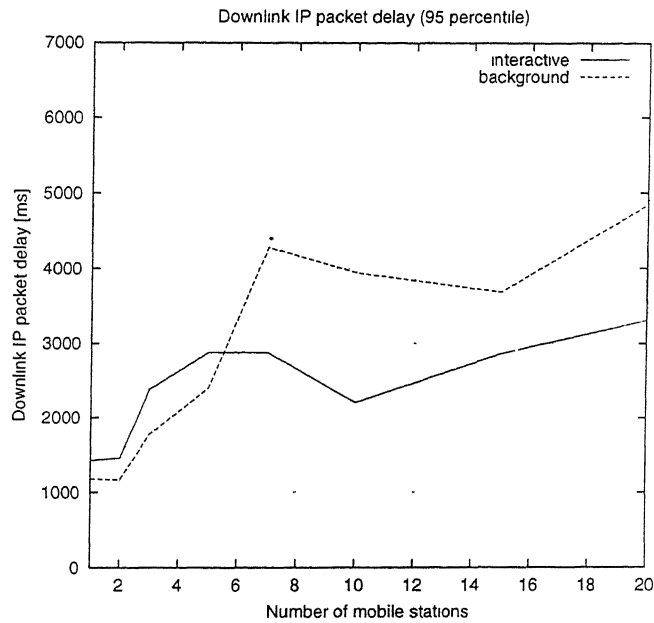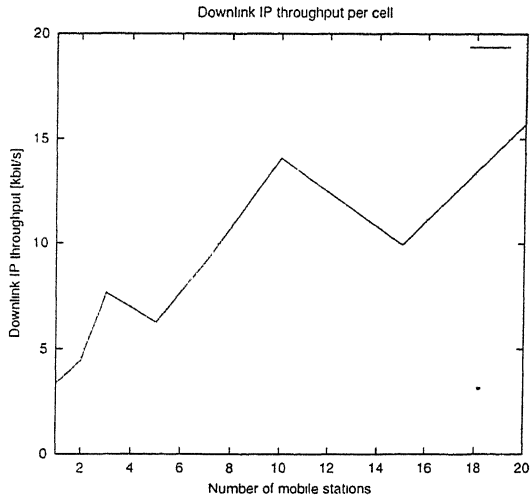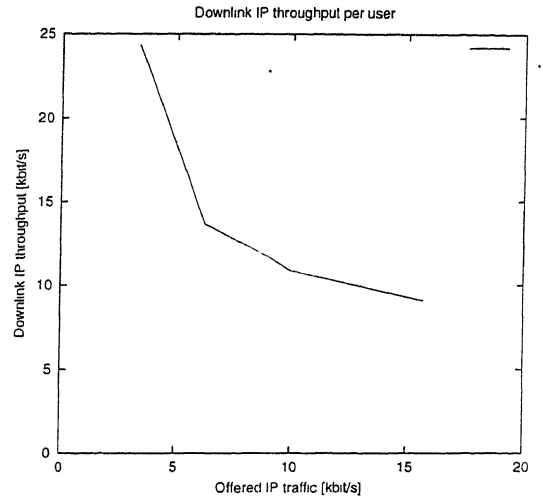
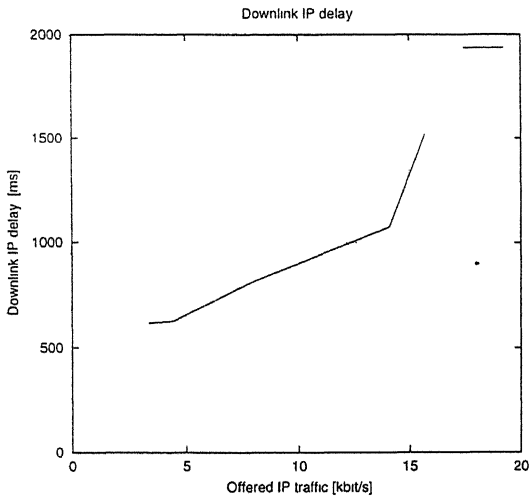Figure 7.5: Downlink IP packet delay (95 percentile)

Further Figure 7.6(a) shows the relationship between downlink IP throughput per cell in Kbps vs number of mobile stations. Followed by Figure 7.6(b) which shows plot of downlink IP throughput per user in Kbps vs offered traffic in Kbps. These graph show that as the number of mobile stations are increased the traffic increases which leads us towards congestion and the downlink IP throughput per user gets decreased. Similarly from Figure 7.6(c) and Figure 7.6(d) it is seen that as traffic in the cell increases the packet delay also increases. Thus it is sure that DiffServ behaves better as traffic is increased because it divides the traffic into different streams and controls quality of service depending upon their priority. Thus from graphs the inherent advantages of differentiated services can be viewed. By carefully enforcing the aggregate traffic profiles and ensuring that new traffic is not admitted that exceeds any aggregate profile, well-defined end-to-end services may be provided over chains of separately administered clouds with DiffServ in the mobile core networks. Furthermore, since each aggregate contracts exists only at the boundary between two clouds, the result is a set of simple bilateral service agreements.
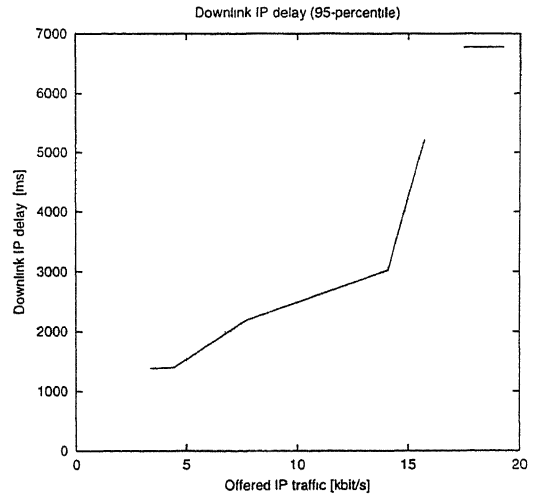
(a) DL IP throughput per cell vs No. of MS



(b) DL IP throughput per user vs Offered traffic



(c) DL IP packet delay vs Offered traffic



(d) DL IP packet delay (95 percentile) vs Offered traffic

Figure 7.6: Other simulation results with DiffServ in mobile core network

## 7.4 Simulation Results with best effort services as compared with DiffServ

In the absence of any specified priority classes and presence of only best effort class the DiffServ IP mobile core network model can be seen as a best effort model. The simulation scenario for the best effort network can be seen from Table 7.2.

Table 7.2: Parameters for simulations with best effort in the IP core network of GPRSim

| Parameter | Value |
|---|---|
| Number of transceiver (TRX) | 3 |
| Model time | One hour |
| Number of MS | 1–20 |
| Packet data channels (PDCHs) fixed | 0 |
| Traffic channels (TCHs) | 21 |
| Packet data channels (PDCHs) on demand | 8 |
| Offered CS traffic [Erlang] | 12.838 |
| $C/I$ ratio [dB] | 12 |
| Session inter-arrival time [s] | 12 |
| Relation *www/email* (Interactive/background sessions) | 3/7 |
| HTTP read time [s] | 12 |
| Pages per session | 5 |
| Objects per page | 2.5 |
| HTTP object size [byte] | 3700 |
| E-mail size [byte] | 10000 |
| Number of standard subscribers [%] | 0 |
| Number of best effort subscribers [%] | 100 |
| IP core output rate | 2Mbps |
| IP core background traffic rate | 1Mbps |
| Buffer size per priority class | 0.1MB |
| schedule time | 40 ms |

Figure 7.7 shows the graph between downlink IP throughput per user in Kbps with respect to number of mobile stations. In this graph the curve for best effort throughput is taken from the above simulation, whereas the curves for interactive and background traffic throughput are taken from Figure 7.1. This is done to compare the best effort services with the DiffServ services. It has been observed that in times when only few mobile stations are active the downlink IP throughput for best effort services is quite good as compared to the downlink IP throughput of interactive and background traffic, but as the number of mobile stations are increased downlink IP throughput per user for best effort services suffers an exponential type
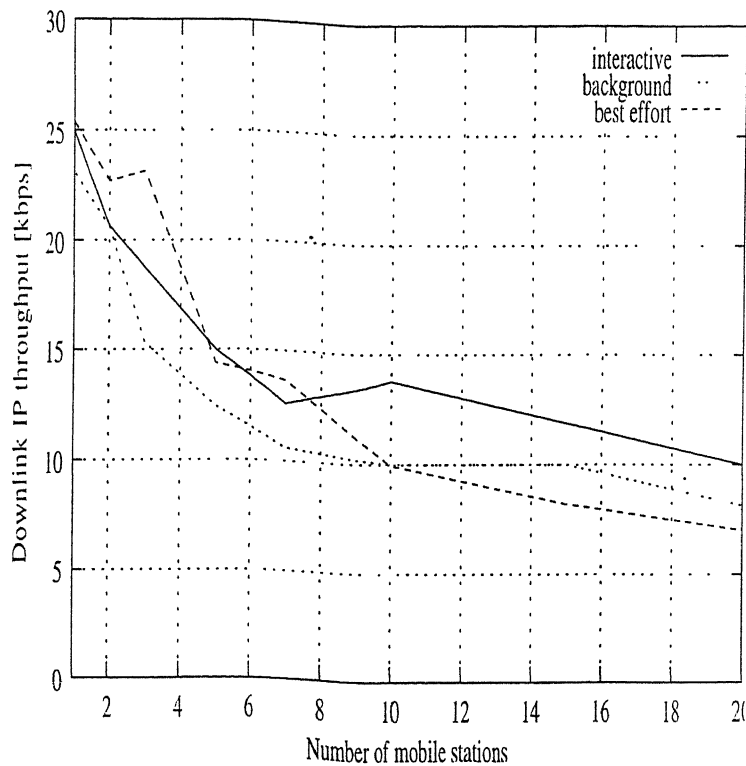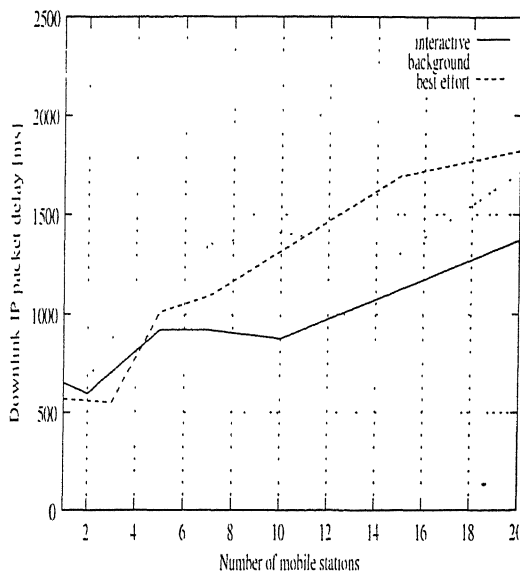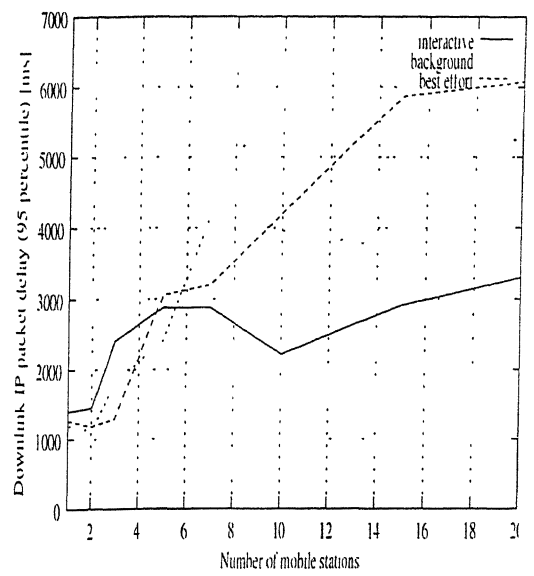
Figure 7.7: Downlink IP throughput with best effort services as compared to DiffServ services

decay which endangers many time sensitive advanced applications. It is seen from Figure 7.7 that as the number of mobile stations become more than 10 the downlink IP throughput of best effort class becomes less than both the background and interactive class. This infers that for high traffic loads DiffServ will be far better than best effort Internet. Thus we see that best effort services will fail to give reliable end-to-end performance guarantee as the traffic on the node increases.

Similarly it can be seen from Figure 7.8(a) and Figure 7.8(b) (which are drawn on the same basis as Figure 7.7) that delay is within bounds for best effort traffic if the number of mobile stations are less but as more and more mobile stations fight for resources the delay (whether it is mean or 95 percentile) increases rapidly making many applications useless. It is seen from these figures that delay of best effort traffic (both mean and 95 percentile) becomes far more as compared to interactive and background class for number of mobile stations more than 12. Thus we see the advantage of using differentiated services for mobile core networks as compared to the present best effort Internet. Hence, it is clear that in cases of high traffic

Figure 7.8: Downlink IP packet delay (mean and 95 percentile) with best effort services as compared to DiffServ services

load conditions DiffServ will steer network out of trouble far better than best effort services.

It is therefore concluded that for IP based mobile core networks, quality of service approach that should be used is differentiated services. It will guarantee under certain conditions non-relative and high performance requirements. Also, it is assumed that there will always exist advanced applications that will stretch the capabilities of the network to the breaking point. Such applications may have bandwidth needs nearly of the order of the aggregate available bottleneck bandwidth, making a reservation free approach impossible. Interoperability, simplicity and scalability are other important strengths of DiffServ approach which makes it even more popular as of today.

# Conclusions

This thesis was intended as a contribution in research towards end-to-end *Quality of Service* deployment in the third generation radio access networks like the *General Packet Radio Service* (GPRS) and the *Universal Mobile Telecommunication System* (UMTS). Specially here the Quality of Service approaches in the developing mobile core networks and the Internet were studied. Enabling quality of service in the third generation core networks is a difficult task as it has to deal with large number of flows at every given instant. The activities that were identified towards reaching the objective of this thesis were:

- Literature study on Internet quality of service approaches like Integrated Services *(IntServ)*, Differentiated Services *(DiffServ)*, and Multi Protocol Label Switching *(MPLS)*.
- Continuous monitoring of the *Third Generation Partnership Project* (3GPP) standardizations for the present and future mobile core networks.
- Continuous monitoring of the latest *Request for comments* (RFCs) and *Internet drafts* of the *Internet Engineering Task Force* (IETF) on differentiated services and other quality of service technologies.
- Designing a general DiffServ architecture as the basis of IP Core Network model (without reference to network configurations and implementation environment).
- Designing in detail the border router architecture residing at the edge of the mobile core network to support interoperability between *3G radio access network* and the *core network.*
- Studying and creating the implementation environment (planning and installing a network configuration, in conjunction with the existing testbed in ComNets)
- Setting up of IP Core Network model based on detailed design of differentiated services architecture.

- Testing the correct behavior of the IP Core Network implementation in presence of ComNets GPRS simulator *GPRSim* to which the IP Core Network model is integrated.

- Comparing the performance gain of DiffServ as compared to the best effort IP for the mobile core networks.

All these activities were fully adhered to and from the results of the previous chapter it has been concluded that *Differentiated Services* (DiffServ) has the capability to support end user quality of service in the *third generation mobile core network* to a great extent. Differentiated services divides the data stream into different priority classes, this makes high priority applications getting superior performance. On the other hand, traffic belonging to less quality of service demanding applications that has been assigned a low priority, does not suffer severe deterioration in quality of service in an order that would have endangered proper application performance and thus user satisfaction.

Thus because of the inherent advantages of differentiated services (i.e. scalable, simple and coarse) it is going to be the technology that will reside in *third generation mobile core networks*. If not alone, it will be there working together with MPLS. Sometimes DiffServ alone cannot assure that the traffic flow will meet absolute bandwidth, delay or jitter requirements as the priority queuing is applied separately at each node and the DiffServ function does not have any knowledge of the load on any other node. MPLS, on the other hand, is a traffic-engineering tool that must take into account the situation at all nodes in the network. The combination of MPLS and DiffServ when implemented in routers will enable the network to create the desired characteristics with suitable guarantees for delay-sensitive traffic and reasonable bandwidth allowances for best-effort traffic.

In this thesis only Differentiated Services have been implemented as at present there is a huge research going on to make IP networks simple and scalable. It is then compared with best effort Internet and the results show clear signs of Differentiated Services advantage.

## 8.1  Future Work

The research to achieve end-to-end quality of service in the Internet and mobile communication will never come to an end, atleast for another five to ten years. Future work may revolve around observing DiffServ vs DiffServ-MPLS combined. Also some other generators should be implemented and integrated into the GPRSIM that allow provisioning of traffic belonging to *Conversational* and *Streaming* classes. More research can be done in the field of IP centric revolution that is attracting the people at present which will make all transfer of user data in form of IP packets. Lastly, quality of service aware Internet service models in mobile IP network environment can be studied and implemented.

# BIBLIOGRAPHY

[AW95]    M.F. Arlitt, C.L. Williamson. *A Synthetic Workload Model for Internet Mosaic Traffic*. In *Proceedings of the 1995 Summer Computer Simulation Conference*. pp. 24 26. Ottawa, Canada, July 1995. Department of Computer Science, University of Sasketchewan, Saskatoon, SK, Canada S7N 5A9.

[BBC+98] S. Blake, D. Black, M. Carlson, E. Davies, Z. Wang, W. Weiss. *An Architecture for Differentiated Services*. Request for Comments 2475, Internet Engineering Task Force. December 1998.

[BBGS00] Y. Bernet, S. Blake, D. Grossman, A. Smith. *An Informal Management Model for DiffServ Routers*. Internet Draft, Internet Engineering Task Force. July 2000.

[BCF00]   S. Brim, B. Carpenter, F. L. Faucher. *Per Hop Behavior Identification Codes*. Request for Comments 2836, Internet Engineering Task Force. May 2000.

[Bla00]    Uyless Black. *QoS in WIDE AREA NETWORKS*. Prentice Hall International. Upper Saddle River, New Jersey, 2000.

[BVE99]   C. Bettstetter, H.J. Vogel, J. Eberspacher. *GSM Phase 2+ General Packet Radio Service (GPRS): Architecture, Protocols and Air Interface*. IEEE Communications Surveys and Tutorials, Vol. 2, No. 3, 1999.

[Cis00]    Cisco Systems. *Cisco Quality of Service Solutions White Paper*. `http://www.cisco.com/warp/public/cc/techno/protocol/tech/qosio.wp.htm`, July 2000.

[DR00]     A. Dutta-Roy. *The cost of quality in Internet-style networks*. IEEE Spectrum. No. 0018-9235/00. pp. 57–62, September 2000.

[EHS97]   J. Ellsberger, D. Hogrefe, A. Sarma. *SDL : Formal Object-oriented Language for Communication Systems*. Prentice Hall Europe, campus 400, Maylands Avenue, Hertfordshire HP2 7EZ, Great Britain, 1997.

[Eri99]    Ericsson ERA/LK. *Ericsson GPRS solutions – End-to-end QoS for GPRS networks*. White Paper ERA/LK-99:00, Ericsson, September 1999. Commercial in Confidence.

[ETS97]  ETSI TR-SMG. *Traffic models (UMTS 50402)*. Technical Report 0.9.3, European Telecommunications Standards Institute, Sophia Antipolis, France, April 1997.

[ETS99]  ETSI 3GPP. *UMTS:QoS Concept and Architecture*. Technical Specification 3G TS 23.107, European Telecommunications Standards Institute. Sophia Antipolis. France, December 1999.

[ETS00a]  ETSI 3GPP. *Universal Mobile Telecommunications System (UMTS); Service aspects; Services and Service Capabilities (3G TS 22.105 version 3.8.0 Release 1999)*. Technical Specification ETSI TS 122 105, European Telecommunications Standards Institute, Sophia Antipolis, France, March 2000.

[FHPW00]  M. Fyro, K. Heikkinen, L. Petersen, P. Wiss. *Media gateway for mobile networks*. Ericsson Review, Vol. 77, pp. 216–223, 2000.

[Fos00]  Carl E. Fossa. *An Analysis of Quality of Service Enhanced Internet Protocols in Third-Generation Wireless Networks*. Project report, Virginia State University. April 2000.

[HFB+99]  J. Heinanen, T. Finland, F. Baker, W. Weiss, J. Wroclawski. *An Assured Forwarding PHB Group*. Request for Comments 2597, Internet Engineering Task Force, June 1999.

[JNP99]  V. Jacobson, K. Nichols, K. Poduri. *An Expedited Forwarding PHB*. Request for Comments 2598, Internet Engineering Task Force, June 1999.

[NBBB98]  K. Nichols, S. Blake, F. Baker, D. Black. *Defination of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 headers*. Request for Comments 2474, Internet Engineering Task Force, December 1998.

[Pax94]  V. Paxson. *Empirically-Derived Analytic Models of Wide-Area TCP Connections*. IEEE/ACM Transactions on Networking, Vol. 2 (4), pp. 316–336, August 1994. ftp://ftp.ee.lbl.gov/papers/WAN-TCP-models.ps.Z.

[PF95]  V. Paxson, S. Floyd. *Wide-Area Traffic: The Failure of Poisson Modeling*. IEEE/ACM Transactions on Networking, Vol. 3, pp. 226–244, June 1995. ftp://ftp.ee.lbl.gov/papers/WAN-poisson.ps.Z.

[Sar00]  B. Sarikaya. *Packet Mode in Wireless Networks: Overview of Transition to Third Generation*. http://www.comsoc.org/tr/tutorial/text/sarikaya.htm, 2000.

[SOL97]   M. Steppler, W. Olzem, C. Lampe. *SDL2CNCL 3.6: SDL-PR to C++ code generator using the ComNets Class Library*. RWTH Aachen, Communication Networks Institute, Kopernikusstr. 16, D-52074 Aachen, August 1997.

[Sta99a]  Stardust Forums Inc. *Frequently Asked Questions about IP Quality of Service*. QoS Forum Draft, September 1999.

[Sta99b]  Stardust Forums Inc. *The Need for QoS*. QoS Forum White Paper, July 1999.

[Sta99c]  Stardust Forums Inc. *QoS Protocols and Architecture*. QoS Forum White Paper. July 1999.

[Stu99]   P. Stuckmann. *Definition of the Number of GSM Traffic Channels Needed to Support Multimedia Applications with the General Packet Radio Service*. Diplomarbeit, RWTH Aachen, Communication Networks Institute, March 1999.

[Tri98]   Trillium Digital Systems, Inc. *Trillium IP Quality of Service White Paper*. http://www.trillium.com/whats-new/wp_ipqos.html, April 1998. Web Version 1072007.12, l.v. 11/2000.

[Tri00]   Trillium Digital Systems, Inc. *Third Generation (3G) Wireless White Paper*. http://www.trillium.com/whats-new/wp_3g.pdf, March 2000.

[Wal99]   B.H. Walke. *Mobile Radio Networks – Networking and Protocols*. John Wiley & Sons, Chichester, April 1999.

[Wap00]   Wapland Inc. *What is GPRS?* Wapland White Paper, February 2000.

[Wit00]   Andreas Witzel. *Control Servers in the core network*. Ericsson Review, Vol. 77, pp. 234–243, 2000.

**133747**

**133747**

**Date Slip**

The book is to be returned on
the date last stamped.